

A Stochastic Majorize-Minimize Subspace Algorithm for Online Penalized Least Squares Estimation

Emilie Chouzenoux and Jean-Christophe Pesquet *

September 27, 2016

Abstract

Stochastic approximation techniques play an important role in solving many problems encountered in machine learning or adaptive signal processing. In these contexts, the statistics of the data are often unknown a priori or their direct computation is too intensive, and they have thus to be estimated online from the observed signals. For batch optimization of an objective function being the sum of a data fidelity term and a penalization (e.g. a sparsity promoting function), Majorize-Minimize (MM) methods have recently attracted much interest since they are fast, highly flexible, and effective in ensuring convergence. The goal of this paper is to show how these methods can be successfully extended to the case when the data fidelity term corresponds to a least squares criterion and the cost function is replaced by a sequence of stochastic approximations of it. In this context, we propose an online version of an MM subspace algorithm and we study its convergence by using suitable probabilistic tools. Simulation results illustrate the good practical performance of the proposed algorithm associated with a memory gradient subspace, when applied to both non-adaptive and adaptive filter identification problems.

Keywords: stochastic approximation, optimization, subspace algorithms, memory gradient methods, descent methods, recursive algorithms, majorization-minimization, filter identification, Newton method, sparsity, machine learning, adaptive filtering.

1 Introduction

A classical problem in data sciences consists of inferring the structure of a linear model linking some observed random variables $(\mathbf{X}_n)_{n \geq 1}$ in $\mathbb{R}^{N \times Q}$ to some other observed random variables $(\mathbf{y}_n)_{n \geq 1}$ in \mathbb{R}^Q . Unless otherwise specified, we will assume in this work that the following wide-sense stationarity properties hold:

$$(\forall n \in \mathbb{N} \setminus \{0\}) \quad \mathbb{E}(\|\mathbf{y}_n\|^2) = \varrho \quad (1)$$

$$\mathbb{E}(\mathbf{X}_n \mathbf{y}_n) = \mathbf{r} \quad (2)$$

$$\mathbb{E}(\mathbf{X}_n \mathbf{X}_n^\top) = \mathbf{R}, \quad (3)$$

*E. Chouzenoux (corresponding author) is with the Laboratoire d'Informatique Gaspard Monge, UMR CNRS 8049, Université Paris-Est, 77454 Marne la Vallée Cedex 2, France. E-mail: emilie.chouzenoux@univ-paris-est.fr. J.-C. Pesquet is with the Center for Visual Computing, CentraleSupélec, University Paris-Saclay, 92295 Chatenay-Malabry, France. E-mail: jean-christophe@pesquet.eu. This work was supported by the CNRS Imag'in project under grant 2015 OPTIMISME. Part of it was presented at the EUSIPCO 2014 conference [1].

where $\varrho \in (0, +\infty)$, $\mathbf{r} \in \mathbb{R}^N$, $\mathbf{R} \in \mathbb{R}^{N \times N}$ is a symmetric positive semi-definite matrix, $\mathbf{E}(\cdot)$ denotes the mathematical expectation, and $\|\cdot\|$ is the Euclidean norm. We will then be interested in the following optimization formulation:

$$\underset{\mathbf{h} \in \mathbb{R}^N}{\text{minimize}} \quad F(\mathbf{h}), \quad (4)$$

with¹

$$(\forall \mathbf{h} \in \mathbb{R}^N) \quad F(\mathbf{h}) = \frac{1}{2} \mathbf{E}(\|\mathbf{y}_n - \mathbf{X}_n^\top \mathbf{h}\|^2) + \Psi(\mathbf{h}), \quad (5)$$

where Ψ is a function from \mathbb{R}^N to \mathbb{R} , playing the role of a regularization function. This penalty function may for instance be useful to incorporate some prior knowledge about the sought parameter vector \mathbf{h} , e.g. some sparsity requirement, possibly in some transformed domain. In this paper, a family of differentiable, non necessarily convex, regularization functions [2] is considered. Problem (4) is encountered in numerous applications such as system identification, channel equalization, linear prediction or interpolation, echo cancellation, interference removal, and supervised classification. In the latter area, $(\mathbf{X}_n)_{n \geq 1}$ are vectors ($Q = 1$) which may correspond to features obtained through some nonlinear mapping of the data to be classified in a given training sequence, and $(\mathbf{y}_n)_{n \geq 1}$ may be the associated (discrete-valued) class index vector [3–5]. Although some other measures (e.g. the logistic regression function) are often more effective in this context, the use of a least squares criterion may still be competitive for simplicity reasons [6, 7], while the regularization term serves here to avoid overfitting which could arise when the number of extracted features is large [8]. Signal reconstruction constitutes another application field of interest. Then, the vector \mathbf{h} corresponds to an unknown signal related to some measurements $(\mathbf{y}_n)_{n \geq 1}$ obtained through products with matrices $(\mathbf{X}_n^\top)_{n \geq 1}$, and additionally corrupted by some noise process [9–11]. Each matrix \mathbf{X}_n^\top with $n \in \mathbb{N} \setminus \{0\}$ corresponds to Q lines of the full acquisition matrix and it is here considered as random. Under suitable stationarity assumptions, the classical least squares data fidelity term can be modeled as $\mathbf{E}(\|\mathbf{y}_n - \mathbf{X}_n^\top \mathbf{h}\|^2)/2$, whereas due to the ill-posedness of the great majority of such inverse problems, a regularization term Ψ needs to be introduced so as to obtain reliable estimates [12].

Many optimization algorithms can be devised to solve Problem (4) depending on the assumptions made on Ψ [13–16]. In this work, we will be interested in Majorize-Minimize (MM) algorithms [17, 18]. In such approaches, the iterates result from successive minimizations of simple surrogates (e.g. quadratic surrogates) majorizing the cost-function. MM algorithms are very flexible and benefit from good theoretical and practical convergence properties. However, the computation load resulting from the minimization of the majorant function may be prohibitive in the context of large scale problems. The strategy we will adopt in this work is to account for subspace acceleration [19], i.e., to constrain the inner minimization step to a subspace of low dimension, typically restricted to the gradient computed at the current iterate and to a memory part (e.g. the difference between the current iterate and a previous one). In a number of recent works [2, 20, 21], MM subspace algorithms provide fast numerical solutions to optimization problems involving smooth functions, in particular in the case of large-scale problems. Note that, although our approach will require that Ψ is a differentiable function, it has been shown that tight approximations of nonsmooth penalizations such as ℓ_1 (resp. ℓ_0) functions, namely $\ell_2 - \ell_1$ (resp. $\ell_2 - \ell_0$) functions, can be employed and are often quite effective in practice [2, 21]. Another advantage of the class of optimization methods under investigation is that their convergence can be established under some technical assumptions, even in the case when Ψ is a nonconvex function (see [2] for more details).

¹The wide sense stationarity assumption makes F independent of the choice of $n \in \mathbb{N} \setminus \{0\}$.

One of the difficulties encountered in machine learning or adaptive processing is that Problem (4) cannot be directly solved since the second-order statistical moments ϱ , \mathbf{r} and \mathbf{R} are often unknown a priori or their direct computation is too intensive, and they have thus to be estimated online. In the simple case when $\Psi = 0$, the classical Recursive Least Squares (RLS) algorithm can be used for this purpose [22]. When Ψ is nonzero, stochastic approximation algorithms have been developed such as the celebrated stochastic gradient descent (SGD) algorithm [23–26] and some of its proximal extensions [27–30]. The convergence speed of SGD may be relatively slow so that various extensions of it have been developed to alleviate this problem (see [9, 31–33] and the references therein). Many efforts have also been devoted to developing adaptive variants of this algorithm [34, 35], in particular when identifying filters having sparse impulse responses (see e.g. [36–42]). In addition, in [43], a set theoretic approach is adopted for online sparse estimation based on projections onto weighted ℓ_1 balls, which is extended in [44] by making use of generalized thresholding mappings. It is worth noting that a sparse RLS algorithm was proposed in [45] for complex-valued signals in the case when Ψ is an ℓ_1 norm. An online variant of the RLS algorithm corresponding to a time weighted LASSO estimator was also designed in [46] which relies on a coordinate descent approach. A similar problem was also addressed in [47] by adopting a novel Bayes variational approach, for which weak theoretical convergence guarantees however exist. If we except [48] where an adaptive primal-dual splitting is employed to deal with a total variation penalization, in almost all the works on sparse adaptive filtering, the sparsity is directly imposed on the filter coefficients, without introducing any linear transform of them.

Designing Majorize-Minimize optimization algorithms in a stochastic context constitutes a challenging task since most of the existing works concerning these methods have been focused on batch optimization procedures, and the related convergence proofs usually rely on deterministic tools. We can however mention a few recent works [49–51] where stochastic MM algorithms have been investigated for general loss functions under specific assumptions (e.g. the independence of the involved random variables [49, 50]), but without introducing any search subspace. Works which are more closely related to ours are those based on Newton or quasi-Newton stochastic algorithms [52–57], in particular the approaches in [54, 55] provide extensions of BFGS algorithm, but proving the convergence of these algorithms requires some specific assumptions. Like BFGS approaches, MM subspace methods use a memory of previous estimates so as to accelerate the convergence.

Our main contributions in this paper are:

- to propose an online version of the MM subspace algorithm from [2, 20], for a wide class of penalized least squares problems,
- to derive a recursive form, with reduced complexity, of the resulting online MM subspace method,
- to prove the convergence of the iterates produced by our method in the stochastic context,
- to show the good practical performance of this method when it is combined with a memory gradient subspace.

In Section 2, we show how Problem (4) can be reformulated in a learning context. The MM strategy which is proposed in this work is described in Section 3.1. In Section 3.3, we give the form of the resulting recursive algorithm and, in Section 3.4, we evaluate its computational complexity. A convergence analysis of the proposed stochastic Majorize-Minimize subspace algorithm is performed in Section 4. In Section 5, two simulation examples in the context of filter identification illustrate the good performance of our algorithm when a memory gradient subspace is employed. Some conclusions are drawn in Section 7.

Table 1: Smooth penalty functions ψ_s fulfilling Assumption 1 and their associated weighting functions ν_s . All expressions are valid for $t \in \mathbb{R}$, $(\lambda_s, \delta_s) \in (0, +\infty)^2$ and $\kappa_s \in [1, 2]$.

| | $\lambda_s^{-1} \psi_s(t)$ | $\lambda_s^{-1} \nu_s(t)$ | Type | Name |
|-----------|--|--|-------------------------------------|--------------------|
| Convex | $ t - \delta_s \log(t /\delta_s + 1)$ | $(t + \delta_s)^{-1}$ | $\ell_2 - \ell_1$ | |
| | $\begin{cases} t^2 & \text{if } t < \delta_s \\ 2\delta_s t - \delta_s^2 & \text{otherwise} \end{cases}$ | $\begin{cases} 2 & \text{if } t < \delta_s \\ 2\delta_s/ t & \text{otherwise} \end{cases}$ | $\ell_2 - \ell_1$ | Huber |
| | $\log(\cosh(t))$ | $\begin{cases} \tanh(t)/t & \text{if } t \neq 0 \\ 1 & \text{otherwise} \end{cases}$ | $\ell_2 - \ell_1$ | Green |
| | $(1 + t^2/\delta_s^2)^{\kappa_s/2} - 1$ | $\kappa_s \delta_s^{-2} (1 + t^2/\delta_s^2)^{\kappa_s/2-1}$ | $\ell_2 - \ell_{\kappa_s}$ | |
| Nonconvex | $\frac{1 - \exp(-t^2/(2\delta_s^2))}{t^2/(2\delta_s^2 + t^2)}$ | $\frac{\delta_s^{-2} \exp(-t^2/(2\delta_s^2))}{4\delta_s^2/(2\delta_s^2 + t^2)}$ | $\ell_2 - \ell_0$ | Welsch |
| | | | $\ell_2 - \ell_0$ | Geman-McClure |
| | $\begin{cases} 1 - (1 - t^2/(6\delta_s^2))^3 & \text{if } t \leq \sqrt{6}\delta_s \\ 1 & \text{otherwise} \end{cases}$ | $\begin{cases} \delta_s^{-2} (1 - t^2/(6\delta_s^2))^2 & \text{if } t \leq \sqrt{6}\delta_s \\ 0 & \text{otherwise} \end{cases}$ | $\ell_2 - \ell_0$ | Tukey biweight |
| | $\tanh(t^2/(2\delta_s^2))$ | $\delta_s^{-2} (\cosh(t^2/(2\delta_s^2)))^{-2}$ | $\ell_2 - \ell_0$ | Hyberbolic tangent |
| | $\log(1 + t^2/\delta_s^2)$ | $2/(t^2 + \delta_s^2)$ | $\ell_2 - \log$ | Cauchy |
| | $1 - \exp(1 - (1 + t^2/(2\delta_s^2))^{\kappa_s/2})$ | $(\kappa_s/(2\delta_s^2))(1 + t^2/(2\delta_s^2))^{\kappa_s/2-1} \exp(1 - (1 + t^2/(2\delta_s^2))^{\kappa_s/2})$ | $\ell_2 - \ell_{\kappa_s} - \ell_0$ | Chouzenoux |

2 Problem formulation

In a learning context, function F can be replaced by a sequence $(F_n)_{n \geq 1}$ of stochastic approximations of it, which are defined as follows: for every $n \in \mathbb{N} \setminus \{0\}$,

$$\begin{aligned}
 (\forall \mathbf{h} \in \mathbb{R}^N) \quad F_n(\mathbf{h}) &= \frac{1}{2\bar{\vartheta}_n} \sum_{k=1}^n \vartheta^{n-k} \|\mathbf{y}_k - \mathbf{X}_k^\top \mathbf{h}\|^2 + \Psi(\mathbf{h}) \\
 &= \frac{1}{2} \rho_n - \mathbf{r}_n^\top \mathbf{h} + \frac{1}{2} \mathbf{h}^\top \mathbf{R}_n \mathbf{h} + \Psi(\mathbf{h}),
 \end{aligned} \tag{6}$$

where $\vartheta \in (0, 1)$,

$$\bar{\vartheta}_n = \sum_{k=0}^{n-1} \vartheta^k = \begin{cases} n & \text{if } \vartheta = 1 \\ \frac{1 - \vartheta^n}{1 - \vartheta} & \text{if } \vartheta \in (0, 1), \end{cases} \tag{7}$$

and ρ_n , \mathbf{r}_n , and \mathbf{R}_n are given by

$$\rho_n = \frac{1}{\bar{\vartheta}_n} \sum_{k=1}^n \vartheta^{n-k} \|\mathbf{y}_k\|^2 \tag{8}$$

$$\mathbf{r}_n = \frac{1}{\bar{\vartheta}_n} \sum_{k=1}^n \vartheta^{n-k} \mathbf{X}_k \mathbf{y}_k \tag{9}$$

$$\mathbf{R}_n = \frac{1}{\bar{\vartheta}_n} \sum_{k=1}^n \vartheta^{n-k} \mathbf{X}_k \mathbf{X}_k^\top. \tag{10}$$

In the case when $\vartheta = 1$, we retrieve the classical sample estimates of ϱ , \mathbf{r} , and \mathbf{R} . When $\vartheta \in (0, 1)$, it can be interpreted as an exponential forgetting factor [22] which may be useful in adaptive processing scenarios (see Section 6).

Hereafter, we will assume that the regularization function Ψ has the following form:

$$(\forall \mathbf{h} \in \mathbb{R}^N) \quad \Psi(\mathbf{h}) = \frac{1}{2} \mathbf{h}^\top \mathbf{V}_0 \mathbf{h} - \mathbf{v}_0^\top \mathbf{h} + \sum_{s=1}^S \psi_s(\|\mathbf{V}_s \mathbf{h} - \mathbf{v}_s\|) \tag{11}$$

where $\mathbf{v}_0 \in \mathbb{R}^N$, $\mathbf{V}_0 \in \mathbb{R}^{N \times N}$ is a symmetric positive semi-definite matrix, and, for every $s \in \{1, \dots, S\}$, $\mathbf{v}_s \in \mathbb{R}^{P_s}$, $\mathbf{V}_s \in \mathbb{R}^{P_s \times N}$, and $\psi_s: \mathbb{R} \rightarrow \mathbb{R}$ is a smooth function. The first two terms in (11) can be viewed as an elastic net penalty [58], while various choices can be made for

the last term. As shown in Table 1, in addition to quadratic regularization functions (obtained when $S = 1$ and $\psi_1 = 0$), $\ell_2 - \ell_1$ functions and smoothed $\ell_2 - \ell_0$ functions constitute standard choices. The matrices $(\mathbf{V}_s)_{1 \leq s \leq S}$ may be set to identity or they may serve to model possible transforms or discrete differentiation operators, and vectors $(\mathbf{v}_s)_{1 \leq s \leq S}$ may be used to define reference values.

Note that the regularization strategy adopted in [46] amounts to replacing Ψ in (6) by $\lambda_n \bar{\Psi}$ where $\bar{\Psi}$ is a (possibly weighted) ℓ_1 norm and $\lambda_n \in [0, +\infty)$. Consistency results can then be established under the assumption that $\vartheta = 1$ and $\lim_{n \rightarrow +\infty} \lambda_n = 0$. Our approach here is different, not only because we are interested in a wide class of regularization functions, but also in the sense that we are looking for a solution to the fully regularized problem (4) instead of a solution to the mean square criterion.

Our objective in the next section will be to propose an efficient recursive method for minimizing functions $(F_n)_{n \geq 1}$.

3 Proposed method

3.1 Majorization property

At each iteration $n \in \mathbb{N} \setminus \{0\}$, we propose to replace F_n by a surrogate function $\Theta_n(\cdot, \mathbf{h}_n)$ based on the current estimate \mathbf{h}_n (computed at the previous iteration). More precisely, a tangent majorant function is chosen such that

$$(\forall \mathbf{h} \in \mathbb{R}^N) \quad F_n(\mathbf{h}) \leq \Theta_n(\mathbf{h}, \mathbf{h}_n) \quad (12)$$

$$F_n(\mathbf{h}_n) = \Theta_n(\mathbf{h}_n, \mathbf{h}_n). \quad (13)$$

For the so-defined MM strategy to be worthwhile, the surrogate function has to be built in such a way that its minimization is simple. For this purpose, the following assumptions will be made on the regularization function Ψ defined in (11):

Assumption 1.

- (i) For every $s \in \{1, \dots, S\}$, ψ_s is an even lower-bounded function, which is continuously differentiable, and $\lim_{t \rightarrow 0, t \neq 0} \dot{\psi}_s(t)/t \in \mathbb{R}$, where $\dot{\psi}_s$ denotes the derivative of ψ_s .
- (ii) For every $s \in \{1, \dots, S\}$, $\psi_s(\sqrt{\cdot})$ is concave on $[0, +\infty)$.
- (iii) There exists $\bar{\nu} \in [0, +\infty)$ such that $(\forall s \in \{1, \dots, S\}) (\forall t \in [0, +\infty)) 0 \leq \nu_s(t) \leq \bar{\nu}$, where $\nu_s(t) = \dot{\psi}_s(t)/t$.²

These assumptions are satisfied by a wide class of functions Ψ [59], in particular those corresponding to the choices of the potential functions $(\psi_s)_{1 \leq s \leq S}$ listed in Table 1.

Assumption 1 implies that each function ψ_s is majorized at every $t \in \mathbb{R}$, by a quadratic function, such that

$$(\forall t' \in \mathbb{R}) \quad \psi_s(t') \leq \psi_s(t) + \dot{\psi}_s(t)(t' - t) + \frac{1}{2} \nu_s(|t|)(t' - t)^2. \quad (14)$$

Note that the above inequality is at the core of iterative reweighted least-squares algorithms [60] and of half quadratic methods [61] for the minimization of penalized quadratic functions. The following majorization then straightforwardly results from (14):

²The function is extended by continuity when $t = 0$.

Proposition 1. *Under Assumption 1, for every $n \in \mathbb{N} \setminus \{0\}$ and $\mathbf{h} \in \mathbb{R}^N$, a tangent majorant of F_n at \mathbf{h} is*

$$\begin{aligned} (\forall \mathbf{h}' \in \mathbb{R}^N) \quad \Theta_n(\mathbf{h}', \mathbf{h}) &= F_n(\mathbf{h}) + \nabla F_n(\mathbf{h})^\top (\mathbf{h}' - \mathbf{h}) \\ &\quad + \frac{1}{2}(\mathbf{h}' - \mathbf{h})^\top \mathbf{A}_n(\mathbf{h})(\mathbf{h}' - \mathbf{h}), \end{aligned} \quad (15)$$

where $\mathbf{A}_n(\mathbf{h})$ is given by

$$\mathbf{A}_n(\mathbf{h}) = \mathbf{R}_n + \mathbf{V}_0 + \mathbf{V}^\top \text{Diag}(\mathbf{b}(\mathbf{h})) \mathbf{V} \in \mathbb{R}^{N \times N} \quad (16)$$

$$\mathbf{V} = [\mathbf{V}_1^\top \dots \mathbf{V}_S^\top]^\top \in \mathbb{R}^{P \times N} \quad (17)$$

$$\mathbf{v} = [\mathbf{v}_1^\top \dots \mathbf{v}_S^\top]^\top \in \mathbb{R}^P \quad (18)$$

with $P = P_1 + \dots + P_S$, and $\mathbf{b}(\mathbf{h}) = (b_i(\mathbf{h}))_{1 \leq i \leq P} \in \mathbb{R}^P$ is such that

$$\mathbf{b}(\mathbf{h}) = \left[\nu_1(\|\mathbf{V}_1 \mathbf{h} - \mathbf{v}_1\|) \mathbf{1}_{P_1}^\top \dots \nu_S(\|\mathbf{V}_S \mathbf{h} - \mathbf{v}_S\|) \mathbf{1}_{P_S}^\top \right]^\top, \quad (19)$$

where $\mathbf{1}_P \in \mathbb{R}^P$ denotes a vector of size P with all entries equal to one.

If, we define, for every $n \in \mathbb{N} \setminus \{0\}$, \mathbf{h}_{n+1} as the minimizer of $\Theta_n(\cdot, \mathbf{h}_n)$, we obtain an online form of a half-quadratic algorithm [61]. Half-quadratic algorithms are known to be effective batch optimization methods, but the use of such method requires the inversion of matrix $\mathbf{A}_n(\mathbf{h}_n)$ at each iteration n , which may be intractable in the context of large scale problems. Subsequently, we propose a subspace acceleration strategy so as to reduce the computational cost of the proposed method.

3.2 Subspace acceleration strategy

The main idea of subspace acceleration is to restrict the minimization space to a subspace spanned by a small number of vectors, instead of minimizing the majorant over the whole space. The proposed MM subspace algorithm consists of defining the following sequence of random vectors $(\mathbf{h}_n)_{n \geq 1}$:

$$(\forall n \in \mathbb{N} \setminus \{0\}) \quad \mathbf{h}_{n+1} \in \arg \min_{\mathbf{h} \in \text{ran } \mathbf{D}_n} \Theta_n(\mathbf{h}, \mathbf{h}_n), \quad (20)$$

where \mathbf{h}_1 has to be set to an initial value, and $\text{ran } \mathbf{D}_n$ denotes the range of a matrix $\mathbf{D}_n \in \mathbb{R}^{N \times M_n}$ that should satisfy the above assumption:

Assumption 2. *For every $n \in \mathbb{N} \setminus \{0\}$, $\{\nabla F_n(\mathbf{h}_n), \mathbf{h}_n\} \subset \text{ran } \mathbf{D}_n$.*

Several approaches can be considered to construct \mathbf{D}_n fulfilling Assumption 2 [20, Tab.I]. The simplest choice is to set $\mathbf{D}_n = [-\nabla F_n(\mathbf{h}_n), \mathbf{h}_n]$, so that (20) reads

$$\mathbf{h}_{n+1} = u_{n,2} \mathbf{h}_n - u_{n,1} \nabla F_n(\mathbf{h}_n), \quad (21)$$

where $(u_{n,1}, u_{n,2})$ is a pair of real-valued random variables. In the special case when $u_{n,2} = 1$, we recover the form of a SGD-like algorithm with step-size $u_{n,1}$. In the machine learning literature, various forms of the step-size for SGD have been proposed [33], which often require to tune up some parameters (e.g. a multiplicative factor) so as to get the best convergence profile on the available dataset. On the contrary, the MM strategy allows us to automatically adjust

$(\mathbf{u}_{n,1}, \mathbf{u}_{n,2})$ at each iteration. Another possibility is to take, for every $n \in \mathbb{N} \setminus \{0\}$, $\text{ran } \mathbf{D}_n = \mathbb{R}^N$. In that case, we recover the online half-quadratic method mentioned earlier, which may have a high computational cost. A more efficient strategy that is at the roots of many works in the context of batch optimization is to adopt an intermediate size subspace matrix, gathering the gradient subspace $[-\nabla F_n(\mathbf{h}_n), \mathbf{h}_n]$ complemented with few vectors containing information regarding the previous iterates (e.g., previous gradient directions, previous iterates,...) [62–64]. In particular, the memory gradient subspace [65], defined as:

$$\mathbf{D}_n = \begin{cases} [-\nabla F_n(\mathbf{h}_n), \mathbf{h}_n, \mathbf{h}_n - \mathbf{h}_{n-1}] & \text{if } n > 1 \\ [-\nabla F_n(\mathbf{h}_1), \mathbf{h}_1] & \text{if } n = 1, \end{cases} \quad (22)$$

was observed to lead to fast convergence on several examples in the field of signal and image restoration [21, 66].

3.3 Recursive MM strategy

We derive in this section a recursive form of the proposed stochastic MM subspace algorithm in (20), with the objective to limit its complexity. First, note that, according to (6), (11), and the definition of functions $(\nu_s)_{1 \leq s \leq S}$ in Assumption 1(iii), for every $n \in \mathbb{N} \setminus \{0\}$, the gradient of F_n is given by

$$(\forall \mathbf{h} \in \mathbb{R}^N) \quad \nabla F_n(\mathbf{h}) = \mathbf{A}_n(\mathbf{h})\mathbf{h} - \mathbf{c}_n(\mathbf{h}), \quad (23)$$

where

$$\mathbf{c}_n(\mathbf{h}) = \mathbf{r}_n + \mathbf{v}_0 + \mathbf{V}^\top \text{Diag}(\mathbf{b}(\mathbf{h}))\mathbf{v} \in \mathbb{R}^N. \quad (24)$$

Thus, using (15), we can rewrite (20) as

$$\mathbf{h}_{n+1} = \mathbf{D}_n \mathbf{u}_n, \quad (25)$$

where \mathbf{u}_n is an \mathbb{R}^{M_n} -valued random vector such that:

$$\begin{aligned} \mathbf{u}_n &= \mathbf{B}_n^\dagger \mathbf{D}_n^\top (\mathbf{A}_n(\mathbf{h}_n)\mathbf{h}_n - \nabla F_n(\mathbf{h}_n)) \\ &= \mathbf{B}_n^\dagger \mathbf{D}_n^\top \mathbf{c}_n(\mathbf{h}_n), \end{aligned} \quad (26)$$

with

$$\mathbf{B}_n = \mathbf{D}_n^\top \mathbf{A}_n(\mathbf{h}_n) \mathbf{D}_n \quad (27)$$

and $(\cdot)^\dagger$ denoting the pseudo-inverse operation. It is important to note that, as \mathbf{B}_n is of dimension $M_n \times M_n$ where M_n is small (typically $M_n = 3$ for the choice of the subspace in (22) when $n > 1$), this pseudo-inversion is light. This constitutes the key advantage of the proposed approach.

By using (7), (9) and (10), the following recursive updates of $(\mathbf{r}_n)_{n \geq 1}$ and $(\mathbf{R}_n)_{n \geq 1}$, can be performed

$$(\forall n \in \mathbb{N} \setminus \{0\}) \quad \mathbf{r}_n = \mathbf{r}_{n-1} + \frac{1}{\vartheta_n} (\mathbf{X}_n \mathbf{y}_n - \mathbf{r}_{n-1}) \quad (28)$$

$$\mathbf{R}_n = \mathbf{R}_{n-1} + \frac{1}{\vartheta_n} (\mathbf{X}_n \mathbf{X}_n^\top - \mathbf{R}_{n-1}), \quad (29)$$

where we have set $\mathbf{r}_0 = \mathbf{0}$ and $\mathbf{R}_0 = \mathbf{O}_N$ and we have used the identity: $\vartheta \bar{\vartheta}_{n-1} / \bar{\vartheta}_n = 1 - \bar{\vartheta}_n^{-1}$. Then, it follows from (16), (27) and (29) that

$$(\forall n \in \mathbb{N} \setminus \{0\}) \quad \mathbf{B}_n = \mathbf{D}_n^\top (\mathbf{D}_n^\mathbf{R} + \mathbf{D}_n^{\mathbf{V}_0}) + (\mathbf{D}_n^\mathbf{V})^\top \text{Diag}(\mathbf{b}(\mathbf{h}_n)) \mathbf{D}_n^\mathbf{V}, \quad (30)$$

where

$$(\forall n \in \mathbb{N} \setminus \{0\}) \quad \mathbf{D}_n^{\mathbf{R}} = \mathbf{R}_n \mathbf{D}_n \in \mathbb{R}^{N \times M_n} \quad (31)$$

$$\mathbf{D}_n^{\mathbf{V}_0} = \mathbf{V}_0 \mathbf{D}_n \in \mathbb{R}^{N \times M_n} \quad (32)$$

$$\mathbf{D}_n^{\mathbf{V}} = \mathbf{V} \mathbf{D}_n \in \mathbb{R}^{P \times M_n}. \quad (33)$$

Finally, let us assume, without loss of generality, that the algorithm is initialized with $\mathbf{h}_1 = \mathbf{D}_0 \mathbf{u}_0$, where $\mathbf{D}_0 \in \mathbb{R}^{N \times M_0}$ and $\mathbf{u}_0 \in \mathbb{R}^{M_0}$. Then, (23) and (25) yield

$$(\forall n \in \mathbb{N} \setminus \{0\}) \quad \nabla F_n(\mathbf{h}_n) = \mathbf{D}_{n-1}^{\mathbf{A}} \mathbf{u}_{n-1} - \mathbf{c}_n(\mathbf{h}_n), \quad (34)$$

where we have set

$$(\forall n \in \mathbb{N}) \quad \mathbf{D}_n^{\mathbf{A}} = \mathbf{A}_{n+1}(\mathbf{h}_{n+1}) \mathbf{D}_n \in \mathbb{R}^{N \times M_n}. \quad (35)$$

By using (16), (29) and (31)-(33), the latter variable can be reexpressed as

$$\begin{aligned} \mathbf{D}_n^{\mathbf{A}} &= \mathbf{R}_{n+1} \mathbf{D}_n + \mathbf{D}_n^{\mathbf{V}_0} + \mathbf{V}^\top \text{Diag}(\mathbf{b}(\mathbf{h}_{n+1})) \mathbf{D}_n^{\mathbf{V}} \\ &= (1 - \frac{1}{\vartheta_{n+1}}) \mathbf{D}_n^{\mathbf{R}} + \frac{1}{\vartheta_{n+1}} \mathbf{X}_{n+1} (\mathbf{X}_{n+1}^\top \mathbf{D}_n) + \mathbf{D}_n^{\mathbf{V}_0} \\ &\quad + \mathbf{V}^\top \text{Diag}(\mathbf{b}(\mathbf{h}_{n+1})) \mathbf{D}_n^{\mathbf{V}}. \end{aligned} \quad (36)$$

The resulting relations are summarized in Algorithm 1, which can be understood as a recursive implementation of Algorithm (20).

Algorithm 1: Stochastic MM subspace method

| | | |
|---|---|---|
| $\mathbf{r}_0 = \mathbf{0}, \mathbf{R}_0 = \mathbf{O}_N$ Initialize $\mathbf{D}_0, \mathbf{u}_0$ $\mathbf{h}_1 = \mathbf{D}_0 \mathbf{u}_0, \mathbf{D}_0^{\mathbf{R}} = \mathbf{O}_{N \times M_n}, \mathbf{D}_0^{\mathbf{V}_0} = \mathbf{V}_0 \mathbf{D}_0, \mathbf{D}_0^{\mathbf{V}} = \mathbf{V} \mathbf{D}_0$ for $n = 1, \dots$ do | $\mathbf{r}_n = \mathbf{r}_{n-1} + \frac{1}{\vartheta_n} (\mathbf{X}_n \mathbf{y}_n - \mathbf{r}_{n-1})$ $\mathbf{c}_n(\mathbf{h}_n) = \mathbf{r}_n + \mathbf{v}_0 + \mathbf{V}^\top \text{Diag}(\mathbf{b}(\mathbf{h}_n)) \mathbf{v}$ $\mathbf{D}_{n-1}^{\mathbf{A}} = (1 - \frac{1}{\vartheta_n}) \mathbf{D}_{n-1}^{\mathbf{R}} + \frac{1}{\vartheta_n} \mathbf{X}_n (\mathbf{X}_n^\top \mathbf{D}_{n-1})$ $\quad + \mathbf{D}_{n-1}^{\mathbf{V}_0} + \mathbf{V}^\top \text{Diag}(\mathbf{b}(\mathbf{h}_n)) \mathbf{D}_{n-1}^{\mathbf{V}}$ $\nabla F_n(\mathbf{h}_n) = \mathbf{D}_{n-1}^{\mathbf{A}} \mathbf{u}_{n-1} - \mathbf{c}_n(\mathbf{h}_n)$ $\mathbf{R}_n = \mathbf{R}_{n-1} + \frac{1}{\vartheta_n} (\mathbf{X}_n \mathbf{X}_n^\top - \mathbf{R}_{n-1})$ Set \mathbf{D}_n using $\nabla F_n(\mathbf{h}_n)$ $\mathbf{D}_n^{\mathbf{R}} = \mathbf{R}_n \mathbf{D}_n, \mathbf{D}_n^{\mathbf{V}_0} = \mathbf{V}_0 \mathbf{D}_n, \mathbf{D}_n^{\mathbf{V}} = \mathbf{V} \mathbf{D}_n$ $\mathbf{B}_n = \mathbf{D}_n^\top (\mathbf{D}_n^{\mathbf{R}} + \mathbf{D}_n^{\mathbf{V}_0}) + (\mathbf{D}_n^{\mathbf{V}})^\top \text{Diag}(\mathbf{b}(\mathbf{h}_n)) \mathbf{D}_n^{\mathbf{V}}$ $\mathbf{u}_n = \mathbf{B}_n^\dagger \mathbf{D}_n^\top (\mathbf{c}_n(\mathbf{h}_n))$ $\mathbf{h}_{n+1} = \mathbf{D}_n \mathbf{u}_n$ | 1 2 3 4 5 6 7 8 9 10 |
| end | | |

3.4 Complexity

Provided that the subspace dimensions $(M_n)_{n \in \mathbb{N}}$ are small, Algorithm 1 has a low complexity, as shown in Table 2.

Table 2: Complexity in terms of multiplications for iteration n of Algorithm 1.

| Step | Complexity for $\mathbf{V} \in \mathbb{R}^{P \times N}$ arbitrary | Complexity when $\mathbf{V} = \mathbf{I}_N$ |
|------|--|--|
| 1 | $N(Q + 1)$ | |
| 2 | $(N + 1)P$ | N |
| 3 | $M_{n-1}(N(2Q + P + 1) + P + Q)$ | $M_{n-1}(N(2Q + 1) + Q)$ |
| 4 | NM_{n-1} | |
| 5 | $N(N + 1)Q/2$ | |
| 7 | $NM_n(2N + P)$ | $2N^2M_n$ |
| 8 | $M_n((M_n + 1)(N + P)/2 + P)$ | $NM_n(M_n + 3)/2$ |
| 9 | $O(M_n^3) + M_n(N + M_n)$ | |
| 10 | NM_n | |

Indeed, the global complexity of a direct implementation of Algorithm 1, evaluated in terms of multiplications at iteration n , is of the order of

$$N(P(M_n + M_{n-1} + 1) + N(4M_n + Q)/2),$$

if we assume that $N \gg \max\{M_n, M_{n-1}, Q\}$. The first term $NP(M_n + M_{n-1} + 1)$ corresponds to an upper bound on the complexity induced by the use of matrices $(\mathbf{V}_s)_{1 \leq s \leq S}$ within the regularization term. Note that these matrices often have a sparse structure (in particular when discrete derivative operators are employed) which may lead to a much lower computational cost. Moreover, when $\mathbf{V} = \mathbf{I}_N$, the identity matrix of \mathbb{R}^N , which is a scenario frequently encountered in adaptive filtering, this term merely vanishes in the evaluation of the global complexity.

The computational complexity can also be reduced by taking advantage of the specific form of matrices $(\mathbf{D}_n)_{n \geq 1}$. Here, we focus our analysis on the example of the memory gradient subspace defined in (22) although it should be noticed that the ideas hereinbelow could be easily generalized to a wide class of subspaces where matrices $(\mathbf{D}_n)_{n \geq 1}$ represent memory features (e.g. [20, Tab. II]). For the particular case of subspace (22), we have:

$$(\forall n > 1) \quad \mathbf{D}_n^{\mathbf{V}} = [-\mathbf{V} \nabla F_n(\mathbf{h}_n), \mathbf{V} \mathbf{h}_n, \mathbf{V} \mathbf{h}_n - \mathbf{V} \mathbf{h}_{n-1}]. \quad (37)$$

Since, for every $n \geq 1$,

$$\mathbf{V} \mathbf{h}_n = \mathbf{V} \mathbf{D}_{n-1} \mathbf{u}_{n-1} = \mathbf{D}_{n-1}^{\mathbf{V}} \mathbf{u}_{n-1}, \quad (38)$$

a recursive formula holds to compute the last two components of $\mathbf{D}_n^{\mathbf{V}}$ in (37). The initial complexity of $3NP$ multiplications is thus reduced to $N(P + 3)$. Similar recursive procedures can be employed to compute $(\mathbf{D}_n^{\mathbf{V}_0})_{n \geq 1}$ allowing the complexity to be reduced to $N(N + 3)$ from $3N^2$. In addition, we have, for every $n > 1$,

$$\mathbf{D}_n^{\mathbf{R}} = [-\mathbf{R}_n \nabla F_n(\mathbf{h}_n), \mathbf{h}_n^{\mathbf{R}}, \mathbf{h}_n^{\mathbf{R}} - \mathbf{R}_n \mathbf{h}_{n-1}], \quad (39)$$

where, by using (29),

$$\begin{aligned} \mathbf{h}_n^{\mathbf{R}} &= \mathbf{R}_n \mathbf{h}_n = (1 - \frac{1}{\vartheta_n}) \mathbf{R}_{n-1} \mathbf{h}_n + \frac{1}{\vartheta_n} \mathbf{X}_n \mathbf{X}_n^{\top} \mathbf{h}_n \\ &= (1 - \frac{1}{\vartheta_n}) \mathbf{D}_{n-1}^{\mathbf{R}} \mathbf{u}_{n-1} + \frac{1}{\vartheta_n} \mathbf{X}_n \mathbf{X}_n^{\top} \mathbf{h}_n \end{aligned} \quad (40)$$

$$\mathbf{R}_n \mathbf{h}_{n-1} = (1 - \frac{1}{\vartheta_n}) \mathbf{h}_{n-1}^{\mathbf{R}} + \frac{1}{\vartheta_n} \mathbf{X}_n \mathbf{X}_n^{\top} \mathbf{h}_{n-1}. \quad (41)$$

It can be further observed that last term $(\bar{\vartheta}_n)^{-1} \mathbf{X}_n \mathbf{X}_n^\top \mathbf{h}_{n-1}$ has already been computed in Step 3 of Algorithm 1. Therefore, instead of $3N^2$ multiplications, we have now to perform $N(N + 2Q + 4)$ ones. With these simplifications, in the case when \mathbf{V}_0 and \mathbf{V} are null matrices, the global complexity of the algorithm is equal to $N^2(Q+2)/2$. When $Q = 1$, we thus recover the order of complexity of the classical RLS algorithm. Since the objective function then reduces to a quadratic function, Sherman-Morrison-Woodbury formula can be invoked to compute iteratively the minimizer on the whole space in an efficient manner.

Note finally that the computation of $\mathbf{X}_n \mathbf{X}_n^\top$ with $n \in \mathbb{N} \setminus \{0\}$, which needs to be performed in Step 5, remains a main source of complexity. However, if $(\forall n > Q) \mathbf{X}_n = [\mathbf{x}_{n-Q+1}, \dots, \mathbf{x}_n]$ where $\mathbf{x}_n \in \mathbb{R}^N$ (as it is the case in affine projection based algorithms for adaptive processing [67]), then a recursive computation of $\mathbf{X}_n \mathbf{X}_n^\top$ only requires $\mathbf{x}_n \mathbf{x}_n^\top$ to be computed at each iteration $n > Q$. If we further assume that the model is a one-dimensional convolutive one, i.e. \mathbf{x}_n corresponds to shifted samples of a signal $(x(n))_{n \geq 1}$, then $(\forall n > N) \mathbf{x}_n = [x(n - N + 1), \dots, x(n)]^\top$ and $\mathbf{x}_n \mathbf{x}_n^\top$ can be itself computed recursively with a complexity of N operations. Such ideas have been deeply investigated in the literature on fast RLS algorithms [68].

4 Convergence study

Establishing the convergence of stochastic approximation algorithms is challenging [23, 29, 69–71]. Throughout this section and the related appendices, it is assumed that $\vartheta = 1$. The underlying probability space being denoted by $(\Omega, \mathcal{F}, \mathbf{P})$, we will say in short that a property is P-a.s. satisfied if this property holds almost surely.

4.1 Assumptions

For every $n \in \mathbb{N} \setminus \{0\}$, let $\mathcal{X}_n = \sigma((\mathbf{X}_k, \mathbf{y}_k)_{1 \leq k \leq n})$ be the sub-sigma algebra of \mathcal{F} generated by $(\mathbf{X}_k, \mathbf{y}_k)_{1 \leq k \leq n}$. In order to give a proof of convergence of the proposed stochastic MM subspace algorithm, we will make the following additional assumption:

Assumption 3.

- (i) $\mathbf{R} + \mathbf{V}_0$ is a positive definite matrix.
- (ii) $((\mathbf{X}_n, \mathbf{y}_n))_{n \geq 1}$ is a stationary ergodic sequence and, for every $n \in \mathbb{N} \setminus \{0\}$, the elements of \mathbf{X}_n and the components of \mathbf{y}_n have finite fourth-order moments.
- (iii) For every $n \in \mathbb{N} \setminus \{0\}$,

$$\mathbb{E}(\|\mathbf{y}_{n+1}\|^2 | \mathcal{X}_n) = \varrho \tag{42}$$

$$\mathbb{E}(\mathbf{X}_{n+1} \mathbf{y}_{n+1} | \mathcal{X}_n) = \mathbf{r} \tag{43}$$

$$\mathbb{E}(\mathbf{X}_{n+1} \mathbf{X}_{n+1}^\top | \mathcal{X}_n) = \mathbf{R}. \tag{44}$$

- (iv) \mathbf{h}_1 is \mathcal{X}_1 -measurable and, for every $n \in \mathbb{N} \setminus \{0\}$, \mathbf{D}_n is \mathcal{X}_n -measurable.

The following asymptotic results will then be useful in the rest of our developments.

Lemma 1. *Under Assumptions 3(ii) and 3(iii), the following properties hold:*

- (i) $(\rho_n)_{n \geq 1}$, $(\mathbf{R}_n)_{n \geq 1}$, and $(\mathbf{r}_n)_{n \geq 1}$ converge P-a.s. to ϱ , \mathbf{R} and \mathbf{r} , respectively

$$\begin{aligned}
(ii) \quad & \sum_{n=1}^{+\infty} n^{-1} |\rho_n - \varrho| < +\infty \quad \text{P-a.s.} \\
& \sum_{n=1}^{+\infty} n^{-1} \|\mathbf{r}_n - \mathbf{r}\| < +\infty \quad \text{P-a.s.} \\
& \sum_{n=1}^{+\infty} n^{-1} \|\mathbf{R}_n - \mathbf{R}\| < +\infty \quad \text{P-a.s.},
\end{aligned}$$

where $\|\cdot\|$ denotes the spectral matrix norm.

Proof. See Appendix A. □

Remark 1.

- (i) Assumptions 3(ii) and 3(iii) are more general than assuming that $((\mathbf{X}_n, \mathbf{y}_n))_{n \geq 1}$ is an independent identically distributed (i.i.d.) sequence and, for every $n \in \mathbb{N} \setminus \{0\}$, the elements of \mathbf{X}_n and the components of \mathbf{y}_n have finite fourth-order moments.
- (ii) Assumption 3(iv) is satisfied as soon as \mathbf{h}_1 is \mathcal{X}_1 -measurable (e.g. \mathbf{h}_1 is deterministic) and the subspace directions, i.e., the columns of \mathbf{D}_n , only depend on $((\mathbf{X}_k, \mathbf{y}_k, \mathbf{h}_k))_{1 \leq k \leq n}$. This is actually the case for the various subspace constructions listed in [20, Tab. I], and, in particular, for the memory gradient subspace given by (22).

4.2 Almost sure convergence

Let us give the following preliminary property:

Lemma 2. Under Assumptions 1, 2 and 3(ii)-3(iii), $(\mathbf{h}_n)_{n \geq 1}$ is P-a.s. bounded.³

Proof. See Appendix B. □

Combining the previous lemma with classical results on the asymptotic behaviour of almost supermartingales, the convergence of the sequence $(F_n(\mathbf{h}_n))_{n \geq 1}$ can be established:

Lemma 3. Under Assumptions 1-3, $(F_n(\mathbf{h}_n))_{n \geq 1}$ is P-a.s. convergent and $((\mathbf{h}_{n+1} - \mathbf{h}_n)^\top \mathbf{A}_n(\mathbf{h}_n)(\mathbf{h}_{n+1} - \mathbf{h}_n))_{n \geq 1}$ is P-a.s. summable.

Proof. See Appendix C. □

Lemma 3 allows us to deduce the following result on the sequence of gradients computed at each iteration of the algorithm:

Lemma 4. Under Assumptions 1-3, $(\|\nabla F_n(\mathbf{h}_n)\|)_{n \geq 1}$ is P-a.s. square-summable.

Proof. See Appendix D. □

By gathering all the previous results, our main convergence results can now be stated:

Proposition 2. Assume that Assumptions 1-3 hold. Then, the following hold:

- (i) The set of cluster points of $(\mathbf{h}_n)_{n \geq 1}$ is almost surely a nonempty compact connected set.

³We say that a sequence of random vectors is almost surely bounded when the norms of all these vectors can be bounded by some random variable with probability 1.

- (ii) Any element of this set is almost surely a critical point of F .
- (iii) If the functions $(\psi_s)_{1 \leq s \leq S}$ are convex, then $(\mathbf{h}_n)_{n \geq 1}$ converges P-a.s. to the unique (global) minimizer of F .

Proof. See Appendix E. □

It can be noticed that the conclusion of Proposition 2(iii) is still valid if the functions $(\psi_s)_{1 \leq s \leq S}$ are nonconvex, they are twice continuously differentiable, and the regularization constants $(\lambda_s)_{1 \leq s \leq S}$ as defined in Table 1 are small enough so that the function F is strongly convex.

4.3 Convergence rate

Based on our recent results in [72], we provide a convergence rate result for Algorithm (20) in the case when the functions $(\psi_s)_{1 \leq s \leq S}$ are convex and twice differentiable.

Proposition 3. *Suppose that Assumptions 1-3 hold. Let $\epsilon \in]0, +\infty[$ be such that $\epsilon \mathbf{I}_N \prec \mathbf{R} + \mathbf{V}_0$. Then, there exists almost surely $n_\epsilon \in \mathbb{N} \setminus \{0\}$ such that, for every $n \geq n_\epsilon$, $\nabla^2 F_n(\mathbf{h}_n) \succeq \mathbf{R} - \epsilon \mathbf{I}_N + \mathbf{V}_0$ and*

$$F_n(\mathbf{h}_{n+1}) - \inf F_n \leq \theta (F_n(\mathbf{h}_n) - \inf F_n) \quad (45)$$

where $\theta \in [0, 1)$.

More details about the expression of the decay rate can be found in [72].

5 Application to 2D system identification

5.1 Problem statement

We first demonstrate the efficiency of the proposed stochastic algorithm in a 2D system identification problem. We consider the following observation model:

$$\mathbf{y} = S(\bar{\mathbf{h}})\mathbf{x} + \mathbf{w}, \quad (46)$$

where $\mathbf{x} \in \mathbb{R}^L$ and $\mathbf{y} \in \mathbb{R}^L$ represent the original and degraded versions of a given image, $\bar{\mathbf{h}} \in \mathbb{R}^N$ is the vectorized version of an unknown two-dimensional blur kernel, S is the linear operator which maps the kernel to its associated Hankel-block Hankel matrix form, and $\mathbf{w} \in \mathbb{R}^L$ represents a realization of an additive noise. When the images \mathbf{x} and \mathbf{y} are of very large size, finding an estimate $\hat{\mathbf{h}} \in \mathbb{R}^N$ of the blur kernel can be quite memory consuming, but one can expect good estimation performance by learning the blur kernel through a sweep of blocks in the dataset.

Let us denote by $\mathbf{X} \in \mathbb{R}^{L \times N}$ the matrix such that $S(\mathbf{h})\mathbf{x} = \mathbf{X}\mathbf{h}$. Then, we propose to define $\hat{\mathbf{h}}$ as a solution to (4), where, for every $n \in \mathbb{N} \setminus \{0\}$, $\mathbf{y}_n \in \mathbb{R}^Q$ and $\mathbf{X}_n^\top \in \mathbb{R}^{Q \times N}$, are subparts of \mathbf{y} and \mathbf{X} , respectively, corresponding to $Q \in \{1, \dots, L\}$ lines of this vector/matrix. For the regularization term Ψ , we consider, for every $s \in \{1, \dots, N\}$ ($S = N$), an isotropic penalization on the gradient between neighboring coefficients of the blur kernel, i.e., $P_s = 2$ and $\mathbf{V}_s = [\Delta_s^h \ \Delta_s^v]^\top$, where $\Delta_s^h \in \mathbb{R}^N$ (resp. $\Delta_s^v \in \mathbb{R}^N$) is the horizontal (resp. vertical) gradient operator applied at pixel s . The smoothness of \mathbf{h} is then enforced by choosing, for every $s \in \{1, \dots, S\}$ and $u \in \mathbb{R}$, $\psi_s(u) = \lambda \sqrt{1 + u^2/\delta^2}$ with $(\lambda, \delta) \in (0, +\infty)^2$. Finally, in order to guarantee the existence of a unique minimizer, the strong convexity of F is imposed by taking $\mathbf{v}_0 = \mathbf{0}$ and $\mathbf{V}_0 = \tau \mathbf{I}_N$, where τ is a small positive value (typically $\tau = 10^{-10}$).

5.2 Simulation results

The original image, presented in Figure 1(a), is a satellite image, of size 4096×4096 pixels. The original blur kernel $\bar{\mathbf{h}}$ with size 21×21 , and the resulting blurred image, which has been corrupted with a zero-mean white Gaussian noise with standard deviation $\sigma = 0.03$ (the blurred signal-to-noise ratio equals 25.7 dB), are displayed in Figures 1(b)(c). Figure 1(d) presents the estimated kernel, using Algorithm 1, with the subspace given by (22), leading to the so-called stochastic MM memory gradient (S3MG) algorithm. Parameters (λ, δ) were adjusted so as to minimize the normalized root mean square estimation error, here equal to 0.064. Figure 2 illustrates the variations of this estimation error with respect to the computation time for the proposed algorithm, the SGD algorithm with a decreasing stepsize proportional to $n^{-1/2}$, the regularized dual averaging (RDA) method with a constant stepsize from [49], and the accelerated stochastic gradient averaging SAGA method with a constant stepsize from [73]. Tests were running on an Intel(R) Xeon(R) E5-2630 @ 2.6GHz using a Matlab 7 implementation. Note that for the latter three algorithms, the stepsize parameter was optimized manually so as to obtain the best performance in terms of convergence speed. Finally, note that all tested algorithms were observed to provide asymptotically the same estimation quality, whatever the size of the blocks. In this example, as illustrated in Figure 3, the best trade-off in terms of convergence speed is obtained for $Q = 256 \times 256$.

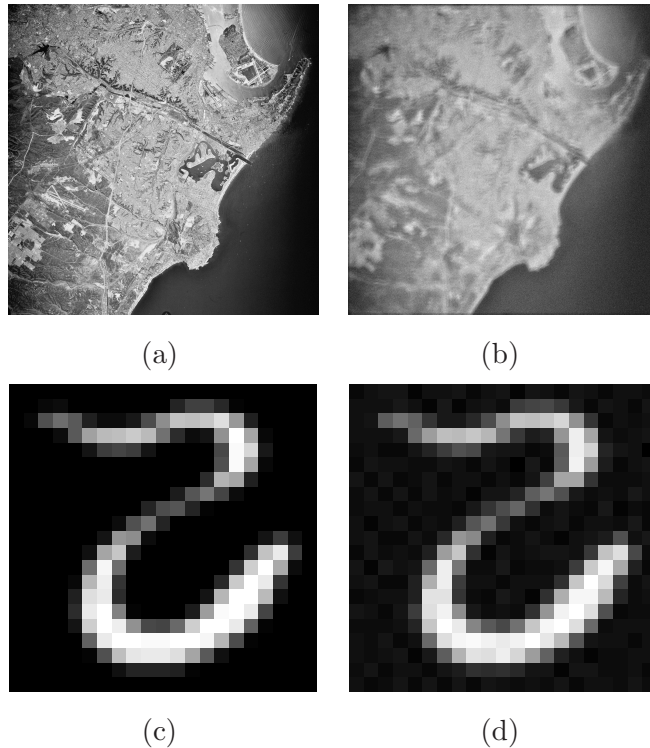


Figure 1: (a) Original image. (b) Blurred and noisy image. (c) Original blur kernel. (d) Estimated blur kernel, with relative error 0.064.

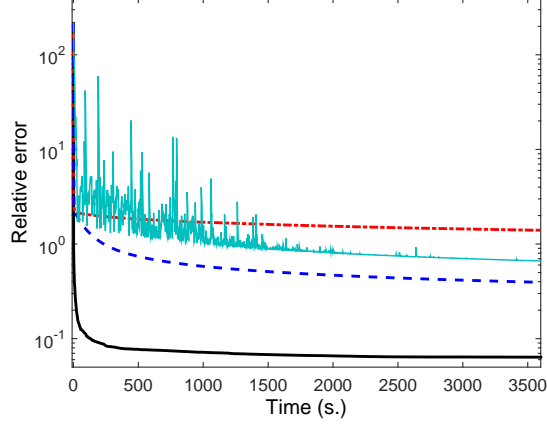


Figure 2: Comparison of S3MG algorithm (solid black line), SGD algorithm with decreasing stepsize $\propto n^{-1/2}$ (dashed-dotted red line), RDA algorithm with constant stepsize (dashed blue line) and SAGA algorithm with constant stepsize (turquoise thin line).

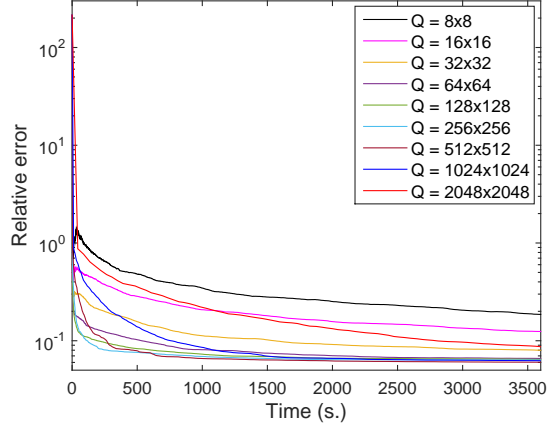


Figure 3: Effect of the block size Q on the convergence speed of S3MG.

6 Application to sparse adaptive filtering

6.1 Problem statement

As emphasized in Sections 2 and 3, one of the advantages of Algorithm 1 compared with some other online optimization algorithms is that it is able to deal with adaptive data processing problems. In this section, we apply the S3MG algorithm to the identification of a sparse time-varying system. Given a real-valued discrete-time input signal $(x(n))_{n \in \mathbb{Z}}$, the output of the system at time $n \geq 1$ is defined as

$$y_n = \mathbf{X}_n^\top \bar{\mathbf{h}}_n + w_n, \quad (47)$$

where $\mathbf{X}_n = [x(n - N + 1), \dots, x(n)]^\top$, w_n models some measurement noise, and $\bar{\mathbf{h}}_n \in \mathbb{R}^N$ gathers the unknown filter taps at time n . Then, the objective is to provide an estimate of the vector $\bar{\mathbf{h}}_n$ at each time by solving Problem (4) where the regularization function Ψ is chosen in order to promote the sparsity of the impulse response of the time-varying filter.

6.2 Simulation results

We generate data according to Model (47) where the input signal $(x(n))_{n \in \mathbb{Z}}$ consists of identically and independent random binary values $\{-1, +1\}$. The measurement noise $(w_n)_{n \in \mathbb{Z}}$ is white Gaussian with zero mean and variance 0.05. In order to evaluate the tracking capability of the proposed S3MG method, the following time-varying linear system is considered:

$$\bar{\mathbf{h}}_n = \begin{cases} \bar{\mathbf{h}}_1 & \text{if } n \leq L/2, \\ \bar{\mathbf{h}}_{L/2+1} & \text{if } n \geq L/2 + 1. \end{cases} \quad (48)$$

The filter length N is equal to 200 and the output of the system is observed at every time $n \in \{1, \dots, L\}$ with $L = 5000$. The sparse impulse responses corresponding to vectors $\bar{\mathbf{h}}_1$ and $\bar{\mathbf{h}}_{L/2+1}$ are represented in Figure 4.

We compute, for every $n \in \{1, \dots, L\}$, the Euclidean norm of the error between the current estimate \mathbf{h}_n and the true filter coefficient vector $\bar{\mathbf{h}}_n$. The minimal estimation error is obtained for the nonconvex Welsch penalty function (see Table 1) and a smoothed $\ell_2 - \ell_0$ regularization function is thus employed by setting $S = N$, $\mathbf{v}_0 = \mathbf{0}$, $\mathbf{V}_0 = \mathbf{O}_N$, and, for every $s \in \{1, \dots, N\}$, $P_s = 1$, $\mathbf{v}_s = \mathbf{0}$, while $\mathbf{V}_s \in \mathbb{R}^{1 \times N}$ is the s -th vector of the canonical basis of \mathbb{R}^N .

We present the results generated by S3MG in Figure 5 for two values of the forgetting factor ϑ , namely $\vartheta = 1$ which corresponds to a non adaptive strategy, and $\vartheta = 0.995$ which appears to be the best choice in terms of tracking properties for this example.

We also show the results obtained with several state-of-the-art approaches in the context of sparse adaptive filtering, namely SPAL [43], RLMS [74], RZAAPA [41] and SM-PAPA [75]. Note that, for each tested method, the involved parameters (stepsize, regularization weight, blocksize) have been tuned manually in order to optimize the performance in terms of error decay.

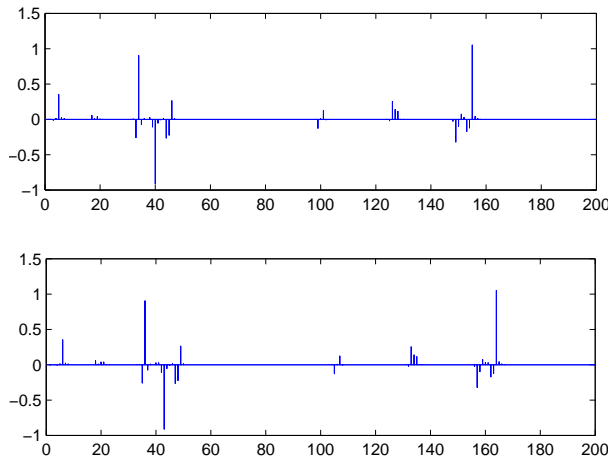


Figure 4: Values of the coefficients of the considered sparse filters $\bar{\mathbf{h}}_1$ (top) and $\bar{\mathbf{h}}_{L/2+1}$ (bottom).

7 Conclusion

In this work, we have proposed a stochastic MM subspace algorithm for online penalized least squares estimation problems. The method makes it possible to use large-size datasets the second-order moments of which are not known a priori. We have shown that the proposed

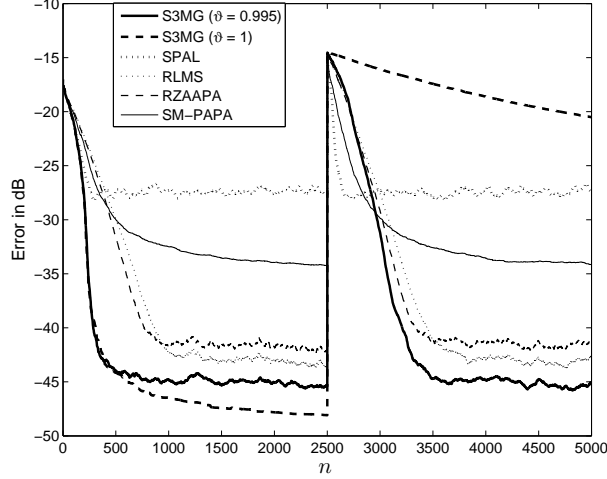


Figure 5: Quadratic estimation error on the filter coefficients as a function of time index n for various adaptive algorithms.

algorithm is of the same order of complexity as the classical RLS algorithm and that its computational cost can be reduced by taking advantage of specific forms of the search subspace. The choice of a memory gradient subspace led to the S3MG algorithm whose good numerical performance has been demonstrated in the context of 2D system identification for large scale image processing problems. In the context of sparse adaptive filtering, S3MG has also been shown to be competitive with respect to recent methods. Although an analysis of the convergence of the proposed method has been carried out, it would be interesting to extend the obtained results to weaker assumptions. In addition, in a nonstationary context, a theoretical study of the tracking abilities of the algorithm should be conducted. Finally, let us emphasize that a detailed analysis of the convergence rate of the proposed method has been undertaken in our recent paper [72].

A Proof of Lemma 1

Property (i) is a consequence of the ergodic theorem [76, Theorem 13.12]. In addition, the law of the iterated logarithm for martingale difference sequences [77] ensures that

$$\limsup_{n \rightarrow +\infty} \frac{|\sum_{k=1}^n (\|\mathbf{y}_k\|^2 - \varrho)|}{(n \log(\log n))^{1/2}} < +\infty \quad \text{P-a.s.} \quad (49)$$

$$\limsup_{n \rightarrow +\infty} \frac{\|\sum_{k=1}^n (\mathbf{X}_k \mathbf{y}_k - \mathbf{r})\|}{(n \log(\log n))^{1/2}} < +\infty \quad \text{P-a.s.} \quad (50)$$

$$\limsup_{n \rightarrow +\infty} \frac{\|\sum_{k=1}^n (\mathbf{X}_k \mathbf{X}_k^\top - \mathbf{R})\|}{(n \log(\log n))^{1/2}} < +\infty \quad \text{P-a.s.} \quad (51)$$

that is

$$\limsup_{n \rightarrow +\infty} \frac{n^{1/2} |\rho_n - \varrho|}{(\log(\log n))^{1/2}} < +\infty \quad \text{P-a.s.} \quad (52)$$

$$\limsup_{n \rightarrow +\infty} \frac{n^{1/2} \|\mathbf{r}_n - \mathbf{r}\|}{(\log(\log n))^{1/2}} < +\infty \quad \text{P-a.s.} \quad (53)$$

$$\limsup_{n \rightarrow +\infty} \frac{n^{1/2} \|\mathbf{R}_n - \mathbf{R}\|}{(\log(\log n))^{1/2}} < +\infty \quad \text{P-a.s.} \quad (54)$$

Consequently, for every $n_0 \in \mathbb{N}$ with $n_0 \geq 2$,

$$\sum_{n=n_0}^{+\infty} n^{-1} |\rho_n - \varrho| \leq \sup_{n \geq n_0} \left(\frac{n^{1/2} |\rho_n - \varrho|}{(\log(\log n))^{1/2}} \right) \left(\sum_{n=n_0}^{+\infty} n^{-3/2} |\log(\log n)|^{1/2} \right). \quad (55)$$

Since $\sum_{n=2}^{+\infty} n^{-3/2} |\log(\log n)|^{1/2} < +\infty$, it follows from (52) that $\sum_{n=n_0}^{+\infty} n^{-1} |\rho_n - \varrho|$ converges P-a.s. to 0 as $n_0 \rightarrow +\infty$, which means that the first line in Property (ii) is satisfied. By proceeding similarly, (53) and (54) allow us to establish the remaining two assertions in Property (ii).

B Proof of Lemma 2

For every $n \in \mathbb{N} \setminus \{0\}$, minimizing $\Theta_n(\cdot, \mathbf{h}_n)$ is equivalent to minimizing the function

$$(\forall \mathbf{h} \in \mathbb{R}^N) \quad \tilde{\Theta}_n(\mathbf{h}, \mathbf{h}_n) = \frac{1}{2} \mathbf{h}^\top \mathbf{A}_n(\mathbf{h}_n) \mathbf{h} - \mathbf{c}_n(\mathbf{h}_n)^\top \mathbf{h}. \quad (56)$$

It follows from Assumption 3(ii)-3(iii) and Lemma 1(i) that there exists $\Lambda \in \mathcal{F}$ such that $P(\Lambda) = 1$ and, for every $\omega \in \Lambda$,

$$\lim_{n \rightarrow +\infty} \mathbf{r}_n(\omega) = \mathbf{r} \quad (57)$$

$$\lim_{n \rightarrow +\infty} \mathbf{R}_n(\omega) = \mathbf{R}. \quad (58)$$

Let $\omega \in \Lambda$. According to Assumption 1(iii) and Eq. (19), $\mathbf{b}(\mathbf{h})$ is bounded as a function of \mathbf{h} . It is then deduced from (24) and (57) that $(\mathbf{c}_n(\mathbf{h}_n)(\omega))_{n \geq 1}$ is bounded, i.e. there exists $\eta \in [0, +\infty)$ such that

$$(\forall n \in \mathbb{N} \setminus \{0\}) \quad \|\mathbf{c}_n(\mathbf{h}_n)(\omega)\| \leq \eta. \quad (59)$$

In addition, as a consequence of (19) and Assumption 1(iii), for every $n \in \mathbb{N} \setminus \{0\}$, $\text{Diag}(\mathbf{b}(\mathbf{h}_n))$ is a positive semidefinite matrix. Hence, because of (16), Assumptions 1(iii) and 3(i), and (58), there exists $\epsilon \in (0, +\infty)$ and $n_0 \in \mathbb{N} \setminus \{0\}$ such that

$$(\forall n \geq n_0) \quad \mathbf{A}_n(\mathbf{h}_n)(\omega) \succeq \mathbf{R} - \epsilon \mathbf{I}_N + \mathbf{V}_0 \succ \mathbf{O}_N. \quad (60)$$

(It suffices to choose ϵ lower than the minimum eigenvalue of $\mathbf{R} + \mathbf{V}_0$). As a consequence of (56), (59), (60), and the Cauchy-Schwarz inequality, we have

$$(\forall n \geq n_0)(\forall \mathbf{h} \in \mathbb{R}^N) \quad \frac{1}{2} \mathbf{h}^\top (\mathbf{R} - \epsilon \mathbf{I}_N + \mathbf{V}_0) \mathbf{h} - \eta \|\mathbf{h}\| \leq \tilde{\Theta}_n(\mathbf{h}, \mathbf{h}_n). \quad (61)$$

Since $\mathbf{R} - \epsilon \mathbf{I}_N + \mathbf{V}_0$ is a positive definite matrix, the lower bound corresponds to a coercive function with respect to \mathbf{h} . There thus exists $\zeta \in (0, +\infty)$ such that, for every $\mathbf{h} \in \mathbb{R}^N$,

$$\|\mathbf{h}\| > \zeta \quad \Rightarrow \quad (\forall n \geq n_0) \quad \tilde{\Theta}_n(\mathbf{h}, \mathbf{h}_n)(\omega) > 0. \quad (62)$$

On the other hand, since $\mathbf{0} \in \text{span}(\mathbf{D}_n(\omega))$, we have

$$\tilde{\Theta}_n(\mathbf{h}_{n+1}, \mathbf{h}_n)(\omega) \leq \tilde{\Theta}_n(\mathbf{0}, \mathbf{h}_n)(\omega) = 0. \quad (63)$$

The last two inequalities allow us to conclude that

$$(\forall n \geq n_0) \quad \|\mathbf{h}_{n+1}(\omega)\| \leq \zeta. \quad (64)$$

C Proof of Lemma 3

According to Assumption 2, the proposed algorithm is actually equivalent to

$$(\forall n \in \mathbb{N} \setminus \{0\}) \quad \mathbf{h}_{n+1} = \mathbf{h}_n + \mathbf{D}_n \tilde{\mathbf{u}}_n \quad (65)$$

$$\tilde{\mathbf{u}}_n = \arg \min_{\tilde{\mathbf{u}} \in \mathbb{R}^M} \Theta_n(\mathbf{h}_n + \mathbf{D}_n \tilde{\mathbf{u}}, \mathbf{h}_n). \quad (66)$$

By using (15) and cancelling the derivative of the function $\tilde{\mathbf{u}} \mapsto \Theta_n(\mathbf{h}_n + \mathbf{D}_n \tilde{\mathbf{u}}, \mathbf{h}_n)$,

$$\mathbf{D}_n^\top \nabla F_n(\mathbf{h}_n) + \mathbf{D}_n^\top \mathbf{A}_n(\mathbf{h}_n) \mathbf{D}_n \tilde{\mathbf{u}}_n = \mathbf{0}. \quad (67)$$

Hence,

$$\begin{aligned} & \Theta(\mathbf{h}_{n+1}, \mathbf{h}_n) \\ &= F_n(\mathbf{h}_n) - \frac{1}{2} \tilde{\mathbf{u}}_n^\top \mathbf{D}_n^\top \mathbf{A}_n(\mathbf{h}_n) \mathbf{D}_n \tilde{\mathbf{u}}_n \\ &= F_n(\mathbf{h}_n) - \frac{1}{2} (\mathbf{h}_{n+1} - \mathbf{h}_n)^\top \mathbf{A}_n(\mathbf{h}_n) (\mathbf{h}_{n+1} - \mathbf{h}_n). \end{aligned} \quad (68)$$

In view of (12) and Proposition 1, this yields

$$(\forall n \in \mathbb{N} \setminus \{0\}) \quad F_n(\mathbf{h}_{n+1}) + \frac{1}{2} (\mathbf{h}_{n+1} - \mathbf{h}_n)^\top \mathbf{A}_n(\mathbf{h}_n) (\mathbf{h}_{n+1} - \mathbf{h}_n) \leq F_n(\mathbf{h}_n). \quad (69)$$

In addition, the following recursive relation holds

$$\begin{aligned} (\forall \mathbf{h} \in \mathbb{R}^N) \quad F_{n+1}(\mathbf{h}) &= F_n(\mathbf{h}) + \frac{1}{2} (\rho_{n+1} - \rho_n) \\ &\quad - (\mathbf{r}_{n+1} - \mathbf{r}_n)^\top \mathbf{h} + \frac{1}{2} \mathbf{h}^\top (\mathbf{R}_{n+1} - \mathbf{R}_n) \mathbf{h}. \end{aligned} \quad (70)$$

As a consequence of Assumption 3(iv), for every $n \in \mathbb{N} \setminus \{0\}$, \mathbf{h}_{n+1} is \mathcal{X}_n -measurable. It can thus be deduced from (69) and the previous two relations that

$$\mathbb{E}(F_{n+1}(\mathbf{h}_{n+1}) | \mathcal{X}_n) + \frac{1}{2} (\mathbf{h}_{n+1} - \mathbf{h}_n)^\top \mathbf{A}_n(\mathbf{h}_n) (\mathbf{h}_{n+1} - \mathbf{h}_n) \leq F_n(\mathbf{h}_n) + \chi_n \quad (71)$$

where

$$\chi_n = \frac{1}{2} \mathbb{E}(\rho_n - \rho_{n+1} | \mathcal{X}_n) - \mathbb{E}(\mathbf{r}_n - \mathbf{r}_{n+1} | \mathcal{X}_n)^\top \mathbf{h}_{n+1} + \frac{1}{2} \mathbf{h}_{n+1}^\top \mathbb{E}(\mathbf{R}_n - \mathbf{R}_{n+1} | \mathcal{X}_n) \mathbf{h}_{n+1}. \quad (72)$$

By using (8)-(10) with $\vartheta = 1$ and Assumption 3(iii), we have

$$\begin{aligned}
\chi_n &= \frac{1}{2(n+1)} (\rho_n - \mathbb{E}(\|\mathbf{y}_{n+1}\|^2 | \mathcal{X}_n)) \\
&\quad - \frac{1}{n+1} (\mathbf{r}_n - \mathbb{E}(\mathbf{X}_{n+1} \mathbf{y}_{n+1} | \mathcal{X}_n))^\top \mathbf{h}_{n+1} \\
&\quad + \frac{1}{2(n+1)} \mathbf{h}_{n+1}^\top (\mathbf{R}_n - \mathbb{E}(\mathbf{X}_{n+1} \mathbf{X}_{n+1}^\top | \mathcal{X}_n)) \mathbf{h}_{n+1} \\
&= \frac{1}{2(n+1)} (\rho_n - \varrho) - \frac{1}{n+1} (\mathbf{r}_n - \mathbf{r})^\top \mathbf{h}_{n+1} \\
&\quad + \frac{1}{2(n+1)} \mathbf{h}_{n+1}^\top (\mathbf{R}_n - \mathbf{R}) \mathbf{h}_{n+1}
\end{aligned} \tag{73}$$

which yields

$$\begin{aligned}
|\chi_n| &\leq \frac{1}{2(n+1)} |\rho_n - \varrho| + \frac{1}{n+1} \|\mathbf{r}_n - \mathbf{r}\| \|\mathbf{h}_{n+1}\| \\
&\quad + \frac{1}{2(n+1)} \|\mathbf{R}_n - \mathbf{R}\| \|\mathbf{h}_{n+1}\|^2.
\end{aligned} \tag{74}$$

According to Lemma 2, $(\mathbf{h}_n)_{n \geq 1}$ is \mathbf{P} -a.s. bounded, and Assumptions 3(ii)-3(iii) and Lemma 1(ii) thus guarantee that

$$\sum_{n=1}^{+\infty} |\chi_n| < +\infty \quad \mathbf{P}\text{-a.s.} \tag{75}$$

Assumption 1(i) entails that, for every $n \in \mathbb{N} \setminus \{0\}$, F_n is lower bounded by $\inf \Psi > -\infty$. Furthermore, (71) leads to

$$\mathbb{E}(F_{n+1}(\mathbf{h}_{n+1}) - \inf \Psi | \mathcal{X}_n) + \frac{1}{2} (\mathbf{h}_{n+1} - \mathbf{h}_n)^\top \mathbf{A}_n(\mathbf{h}_n) (\mathbf{h}_{n+1} - \mathbf{h}_n) \leq F_n(\mathbf{h}_n) - \inf \Psi + |\chi_n|. \tag{76}$$

Since, for every $n \in \mathbb{N} \setminus \{0\}$, $F_n(\mathbf{h}_n) - \inf \Psi$ and $(\mathbf{h}_{n+1} - \mathbf{h}_n)^\top \mathbf{A}_n(\mathbf{h}_n) (\mathbf{h}_{n+1} - \mathbf{h}_n)$ are nonnegative, $(F_n(\mathbf{h}_n) - \inf \Psi)_{n \geq 1}$ is a nonnegative almost supermartingale [78]. By invoking now Siegmund-Robbins lemma [79], it can be deduced from (75) that the desired convergence results hold.

D Proof of Lemma 4

According to (15), we have, for every $\phi \in \mathbb{R}$ and $n \in \mathbb{N} \setminus \{0\}$,

$$\Theta_n(\mathbf{h}_n - \phi \nabla F_n(\mathbf{h}_n), \mathbf{h}_n) = F_n(\mathbf{h}_n) - \phi \|\nabla F_n(\mathbf{h}_n)\|^2 + \frac{\phi^2}{2} (\nabla F_n(\mathbf{h}_n))^\top \mathbf{A}_n(\mathbf{h}_n) \nabla F_n(\mathbf{h}_n). \tag{77}$$

Let

$$\Phi_n \in \underset{\phi \in \mathbb{R}}{\text{Argmin}} \Theta_n(\mathbf{h}_n - \phi \nabla F_n(\mathbf{h}_n), \mathbf{h}_n). \tag{78}$$

The following optimality condition holds:

$$(\nabla F_n(\mathbf{h}_n))^\top \mathbf{A}_n(\mathbf{h}_n) \nabla F_n(\mathbf{h}_n) \Phi_n = \|\nabla F_n(\mathbf{h}_n)\|^2. \tag{79}$$

As a consequence of Assumption 2, $(\forall \phi \in \mathbb{R}) \mathbf{h}_n - \phi \nabla F_n(\mathbf{h}_n) \in \text{span } \mathbf{D}_n$. It then follows from (20) and (79) that

$$\begin{aligned}
\Theta_n(\mathbf{h}_{n+1}, \mathbf{h}_n) &\leq \Theta_n(\mathbf{h}_n - \Phi_n \nabla F_n(\mathbf{h}_n), \mathbf{h}_n) \\
&\leq F_n(\mathbf{h}_n) - \frac{\Phi_n^2}{2} \|\nabla F_n(\mathbf{h}_n)\|^2
\end{aligned} \tag{80}$$

which, by using (68), leads to

$$\Phi_n \|\nabla F_n(\mathbf{h}_n)\|^2 \leq (\mathbf{h}_{n+1} - \mathbf{h}_n)^\top \mathbf{A}_n(\mathbf{h}_n)(\mathbf{h}_{n+1} - \mathbf{h}_n). \quad (81)$$

Let $\epsilon > 0$. Assumption 1(iii) and (16) yield, for every $n \in \mathbb{N} \setminus \{0\}$,

$$\mathbf{A}_n(\mathbf{h}_n) \preceq (\|\mathbf{R}_n + \mathbf{V}_0\| + \bar{\nu} \|\mathbf{V}\|^2) \mathbf{I}_N. \quad (82)$$

Therefore, according to Assumptions 3(i) and 3(ii), and Lemma 1(i), there exists $\Lambda \in \mathcal{F}$ such that $\mathbf{P}(\Lambda) = 1$ and, for every $\omega \in \Lambda$,

$$(\exists n_0 \in \mathbb{N} \setminus \{0\})(\forall n \geq n_0) \quad \mathbf{O}_N \prec \mathbf{A}_n(\mathbf{h}_n)(\omega) \preceq \alpha_\epsilon^{-1} \mathbf{I}_N \quad (83)$$

where

$$\alpha_\epsilon = (\|\mathbf{R} + \mathbf{V}_0\| + \bar{\nu} \|\mathbf{V}\|^2 + \epsilon)^{-1} > 0. \quad (84)$$

Let $\omega \in \Lambda$. By using now (79), it can be deduced from (83) that, if $n \geq n_0$ and $\nabla F_n(\mathbf{h}_n)(\omega) \neq \mathbf{0}$, then

$$\Phi_n(\omega) \geq \alpha_\epsilon. \quad (85)$$

Then, it follows from (81) that

$$\begin{aligned} & \alpha_\epsilon \sum_{n=n_0}^{+\infty} \|\nabla F_n(\mathbf{h}_n)(\omega)\|^2 \\ & \leq \sum_{n=n_0}^{+\infty} (\mathbf{h}_{n+1}(\omega) - \mathbf{h}_n(\omega))^\top \mathbf{A}_n(\mathbf{h}_n)(\omega) (\mathbf{h}_{n+1}(\omega) - \mathbf{h}_n(\omega)). \end{aligned} \quad (86)$$

By invoking Lemma 3, we can conclude that $(\|\nabla F_n(\mathbf{h}_n)\|^2)_{n \geq 1}$ is P-a.s. summable.

E Proof of Proposition 2

It follows from Lemma 3 that $((\mathbf{h}_{n+1} - \mathbf{h}_n)^\top \mathbf{A}_n(\mathbf{h}_n)(\mathbf{h}_{n+1} - \mathbf{h}_n))_{n \geq 1}$ converges P-a.s. to 0. In addition, we have seen in the proof of Lemma 2 that there exists $\Lambda \in \mathcal{F}$ such that $\mathbf{P}(\Lambda) = 1$ and, for every $\omega \in \Lambda$, (60) holds with $\epsilon \in (0, +\infty)$ and $n_0 \in \mathbb{N} \setminus \{0\}$. This implies that, for every $n \geq n_0$,

$$\begin{aligned} & \|\mathbf{R} - \epsilon \mathbf{I}_N + \mathbf{V}_0\| \|\mathbf{h}_{n+1}(\omega) - \mathbf{h}_n(\omega)\|^2 \\ & \leq (\mathbf{h}_{n+1}(\omega) - \mathbf{h}_n(\omega))^\top \mathbf{A}_n(\mathbf{h}_n)(\omega) (\mathbf{h}_{n+1}(\omega) - \mathbf{h}_n(\omega)) \end{aligned} \quad (87)$$

where $\|\mathbf{R} - \epsilon \mathbf{I}_N + \mathbf{V}_0\| > 0$. Consequently, $(\mathbf{h}_{n+1} - \mathbf{h}_n)_{n \geq 1}$ converges P-a.s. to $\mathbf{0}$. In addition, according to Lemma 2, $(\mathbf{h}_n)_{n \geq 1}$ belongs almost surely to a compact set. The result is then obtained by invoking Ostrowski's theorem [80, Theorem 26.1].

(ii) By using (23)-(24), we have

$$(\forall n \in \mathbb{N} \setminus \{0\}) \quad \nabla F_n(\mathbf{h}_n) - \nabla F(\mathbf{h}_n) = (\mathbf{R}_n - \mathbf{R})\mathbf{h}_n - \mathbf{r}_n + \mathbf{r}. \quad (88)$$

Since $(\mathbf{h}_n)_{n \geq 1}$ is almost surely bounded, it follows from Lemma 1(i) that $(\nabla F_n(\mathbf{h}_n) - \nabla F(\mathbf{h}_n))_{n \geq 1}$ converges P-a.s. to $\mathbf{0}$. Since Lemma 4 ensures that $(\nabla F_n(\mathbf{h}_n))_{n \geq 1}$ converges P-a.s. to $\mathbf{0}$, $(\nabla F(\mathbf{h}_n))_{n \geq 1}$ also converges P-a.s. to $\mathbf{0}$. There thus exists $\Lambda \in \mathcal{F}$ such that $\mathbf{P}(\Lambda) = 1$ and, for every $\omega \in \Lambda$, $\nabla F(\mathbf{h}_n(\omega)) \rightarrow \mathbf{0}$. Let $\hat{\mathbf{h}}$ be a cluster point of $(\mathbf{h}_n(\omega))_{n \geq 1}$. There exists a

subsequence $(\mathbf{h}_{k_n}(\omega))_{n \geq 1}$ such that $\mathbf{h}_{k_n}(\omega) \rightarrow \hat{\mathbf{h}}$. As we have assumed that the regularization functions $(\psi_s)_{1 \leq s \leq S}$ are continuously differentiable (see Assumption 1(i)), F is also continuously differentiable, and

$$\nabla F(\hat{\mathbf{h}}) = \lim_{n \rightarrow +\infty} \nabla F(\mathbf{h}_{k_n}(\omega)) = \mathbf{0}. \quad (89)$$

This means that $\hat{\mathbf{h}}$ is a critical point of F .

(iii) Because of Assumption 3(i), when the functions $(\psi_s)_{1 \leq s \leq S}$ are convex, F is a strongly convex function. It thus possesses a unique critical point $\hat{\mathbf{h}}$, which is the global minimizer of F . It follows from (i) and (ii) that, almost surely, the unique cluster point of $(\mathbf{h}_n)_{n \geq 1}$ is $\hat{\mathbf{h}}$, which shows that $\mathbf{h}_n \rightarrow \hat{\mathbf{h}}$ P-a.s.

Acknowledgements

The authors would like to thank Professors Markus V. S. Lima and Paulo S. R. Diniz from Federal University of Rio de Janeiro for kindly making us accessible some of their codes. We would also like to thank Professor Anisia Florescu from Dunărea de Jos University of Galați for providing the initial motivation for this work.

References

- [1] E. Chouzenoux, J.-C. Pesquet, and A. Florescu, “A stochastic 3MG algorithm with application to 2D filter identification,” in *Proc. 22nd European Signal Process. Conf. (EUSIPCO 2014)*, Lisboa, Portugal, 1-5 Sep. 2014, pp. 1587–1591.
- [2] E. Chouzenoux, A. Jezierska, J.-C. Pesquet, and H. Talbot, “A majorize-minimize subspace approach for ℓ_2 - ℓ_0 image regularization,” *SIAM J. Imag. Sci.*, vol. 6, no. 1, pp. 563–591, 2013.
- [3] L. Bottou, “Stochastic learning,” in *Advanced Lectures on Machine Learning*, O. Bousquet and U. von Luxburg, Eds., Lecture Notes in Artificial Intelligence, LNAI 3176, pp. 146–168. Springer Verlag, Berlin, 2004.
- [4] S. Sra, S. Nowozin, and S. J. Wright (eds.), *Optimization for Machine Learning*, MIT Press, Cambridge, MA, 2011.
- [5] S. Theodoridis, *Machine Learning: A Bayesian and Optimization Perspective*, Academic Press, San Diego, CA, 2015.
- [6] R. Rifkin, G. Yeo, and T. Poggio, vol. 190 of *III-Computer and Systems Sciences*, chapter 7, pp. 131–154, 2003.
- [7] A. Tacchetti, P. Mallapragada, M. Santoro, and L. Rosasco, “GURLS: a least squares library for supervised learning,” *J. Mach. Learn. Res.*, vol. 14, pp. 3201–3205, 2013.
- [8] F. Bauer, S. Pereverzev, and L. Rosasco, “On regularization algorithms in learning theory,” *J. Complexity*, vol. 23, no. 1, pp. 52–72, 2007.
- [9] M. Pereyra, P. Schniter, E. Chouzenoux, J.-C. Pesquet, J.-Y. Tournieret, A. O Hero, and S. McLaughlin, “A survey of stochastic simulation and optimization methods in signal processing,” *IEEE J. Sel. Top. Signal Process.*, vol. 10, no. 2, pp. 224–241, Mar. 2016.

- [10] N. Pustelnik, A. Benazza-Benhayia, Y. Zheng, and J.-C. Pesquet, “Wavelet-based image deconvolution and reconstruction,” 2016, to appear in Wiley Encyclopedia of Electrical and Electronics Engineering.
- [11] P. S. R. Diniz, *Adaptive Filtering. Algorithms and Practical Implementation*, Springer, New York, NY, 4th edition, 2013.
- [12] G. Demoment, “Image reconstruction and restoration: Overview of common estimation structures and problems,” *IEEE Trans. Acous., Speech Signal Process.*, vol. 37, no. 12, pp. 2024–2036, 1989.
- [13] J. Nocedal and S. J. Wright, *Numerical Optimization*, Springer-Verlag, New York, 1999.
- [14] P. L. Combettes and J.-C. Pesquet, “Proximal splitting methods in signal processing,” in *Fixed-Point Algorithms for Inverse Problems in Science and Engineering*, H. H. Bauschke, R. Burachik, P. L. Combettes, V. Elser, D. R. Luke, and H. Wolkowicz, Eds., pp. 185–212. Springer-Verlag, New York, 2010.
- [15] M. V. Afonso, J. M. Bioucas-Dias, and M. A. T. Figueiredo, “An augmented Lagrangian approach to the constrained optimization formulation of imaging inverse problems,” *IEEE Trans. Image Process.*, vol. 20, no. 3, pp. 681–695, 2011.
- [16] S. Boyd, N. Parikh, E. Chu, B. Peleato, and J. Eckstein, “Distributed optimization and statistical learning via the alternating direction method of multipliers,” *Found. Trends Machine Learn.*, vol. 8, no. 1, pp. 1–122, 2011.
- [17] D. R. Hunter and K. Lange, “A tutorial on MM algorithms,” *Amer. Stat.*, vol. 58, no. 1, pp. 30–37, Feb. 2004.
- [18] Z. Zhang, J. T. Kwok, and D.-Y. Yeung, “Surrogate maximization/minimization algorithms and extensions,” *Mach. Learn.*, vol. 69, pp. 1–33, 2007.
- [19] M. Zibulevsky and M. Elad, “ $\ell_2 - \ell_1$ optimization in signal and image processing,” *IEEE Signal Process. Mag.*, vol. 27, pp. 76–88, May 2010.
- [20] E. Chouzenoux, J. Idier, and S. Moussaoui, “A majorize-minimize subspace strategy for subspace optimization applied to image restoration,” *IEEE Trans. Image Process.*, vol. 20, no. 18, pp. 1517–1528, Jun. 2011.
- [21] A. Florescu, E. Chouzenoux, J.-C. Pesquet, P. Ciuciu, and S. Ciochina, “A majorize-minimize memory gradient method for complex-valued inverse problem,” *Signal Process.*, vol. 103, pp. 285–295, Oct. 2014, Special issue on Image Restoration and Enhancement: Recent Advances and Applications.
- [22] S. O. Haykin, *Adaptive Filter Theory*, Prentice Hall, New Jersey, USA, 4th edition, 2002.
- [23] H. Robbins and S. Monro, “A stochastic approximation method,” *Ann. Math. Stat.*, vol. 22, no. 3, pp. 400–407, 1951.
- [24] J. M. Ermoliev and Z. V. Nekrylova, “The method of stochastic gradients and its application,” in *Seminar: Theory of Optimal Solutions. No. 1 (Russian)*, pp. 24–47. Akad. Nauk Ukrain. SSR, Kiev, 1967.
- [25] O. V. Guseva, “The rate of convergence of the method of generalized stochastic gradients,” *Kibernetika (Kiev)*, , no. 4, pp. 143–145, 1971.

- [26] D. P. Bertsekas and J. N. Tsitsiklis, “Gradient convergence in gradient methods with errors,” *SIAM J. Optim.*, vol. 10, no. 3, pp. 627–642, 2000.
- [27] Y. F. Atchadé, G. Fort, and E. Moulines, “On stochastic proximal gradient algorithms,” Tech. Rep., 2014, <http://arxiv.org/abs/1402.2365>.
- [28] L. Rosasco, S. Villa, and B. C. Vũ, “Convergence of stochastic proximal gradient algorithm,” Tech. Rep., 2014, <http://arxiv.org/abs/1403.5074>.
- [29] P. L. Combettes and J.-C. Pesquet, “Stochastic approximations and perturbations in forward-backward splitting for monotone operators,” *Pure Appl. Funct. Anal.*, vol. 1, no. 1, pp. 13–37, Jan. 2016.
- [30] J. Konečný, J. Liu, P. Richtárik, and M. Takáč, “Mini-batch semi-stochastic gradient descent in the proximal setting,” *IEEE J. Sel. Top. Signal Process.*, vol. 10, no. 2, pp. 242–255, Mar. 2016.
- [31] B. T. Polyak and A. B. Juditsky, “Acceleration of stochastic approximation by averaging,” *SIAM J. Control Optim.*, vol. 30, no. 4, pp. 838–855, 1992.
- [32] A. Nemirovski, A. Juditsky, G. Lan, and A. Shapiro, “Robust stochastic approximation approach to stochastic programming,” *SIAM J. Optim.*, vol. 19, no. 4, pp. 1574–1609, 2008.
- [33] F. Bach and E. Moulines, “Non-asymptotic analysis of stochastic approximation algorithms for machine learning,” in *Proc. Ann. Conf. Neur. Inform. Proc. Syst.*, Granada, Spain, Dec. 12 - 17 2011, pp. x–x+8.
- [34] B. Widrow and S.D. Stearns, *Adaptive Signal Processing*, Prentice Hall, New Jersey, 1985.
- [35] O. Macchi, *Adaptive Processing: The Least Mean Squares Approach With Applications in Transmission*, Wiley, Chichester, UK, 1995.
- [36] O. Hoshuyama, R. A. Goubran, and A. Sugiyama, “A generalized proportionate variable step-size algorithm for fast changing acoustic environments,” in *Proc. Int. Conf. Acoust., Speech Signal Process. (ICASSP 2004)*, Montreal, Canada, May 17-21 2004, vol. 4, pp. 161–164.
- [37] A. W. H. Khong and P. A. Naylor, “Efficient use of sparse adaptive filters,” in *Proc. Asilomar Conf. Signal, Systems and Computers*, Pacific Grove, CA, Oct. 29-Nov. 1 2006, pp. 1375–1379.
- [38] Y. Chen, Y. Gu, and A. O. Hero, “Sparse LMS for system identification,” in *Proc. Int. Conf. Acoust., Speech Signal Process. (ICASSP 2009)*, Taipei, Taiwan, Apr. 19-24 2009, pp. 3125–3128.
- [39] C. Paleologu, J. Benesty, and S. Ciochină, *Sparse adaptive filters for echo cancellation*, Synthesis Lectures on Speech and Audio Processing. Morgan and Claypool, San Rafael, USA, 2010.
- [40] Y. Murakami, M. Yamagishi, M. Yukawa, and I. Yamada, “A sparse adaptive filtering using time-varying soft-thresholding techniques,” in *Proc. Int. Conf. Acoust., Speech Signal Process. (ICASSP 2010)*, Dallas, Texas, Mar. 14-19 2010, pp. 3734–3737.

- [41] R. Meng, R. C. De Lamare, and V. H. Nascimento, "Sparsity-aware affine projection adaptive algorithms for system identification," in *Proc. Sensor Signal Process. Defence*, London, U.K., Sept. 27-29 2011, pp. 1–5.
- [42] M. V. S. Lima, W. A. Martins, and P. S. R. Diniz, "Affine projection algorithms for sparse system identification," in *Proc. Int. Conf. Acoust., Speech Signal Process. (ICASSP 2013)*, Vancouver, Canada, May 26-31 2013, pp. 5666–5670.
- [43] Y. Kopsinis, K. Slavakis, and S. Theodoridis, "Online sparse system identification and signal reconstruction using projections onto weighted ℓ_1 balls," *IEEE Trans. Signal Process.*, vol. 59, no. 3, pp. 936–952, Mar. 2011.
- [44] K. Slavakis, Y. Kopsinis, S. Theodoridis, and S. McLaughlin, "Generalized thresholding and online sparsity-aware learning in a union of subspaces," *IEEE Trans. Signal Process.*, vol. 61, no. 15, pp. 3760–3773, Aug. 2013.
- [45] B. Babadi, N. Kalouptsidis, and V. Tarokh, "SPARLS: The sparse RLS algorithm," *IEEE Trans. Signal Process.*, vol. 58, no. 8, pp. 4013–4025, Aug. 2010.
- [46] D. Angelosante, J. A. Bazerque, and G. B. Giannakis, "Online adaptive estimation of sparse signals: where RLS meets the ℓ_1 -norm," *IEEE Trans. Signal Process.*, vol. 58, no. 7, pp. 3436–3447, Jul. 2010.
- [47] K. E. Themelis, A. A. Rontogiannis, and K. D. Koutroumbas, "A variational Bayes framework for sparse adaptive estimation," *IEEE Trans. Signal Process.*, vol. 62, no. 18, pp. 4723–4736, 2014.
- [48] S. Ono, M. Yamagishi, and I. Yamada, "A sparse system identification by using adaptively-weighted total variation via a primal-dual splitting approach," in *Proc. Int. Conf. Acoust., Speech Signal Process. (ICASSP 2013)*, Vancouver, Canada, 26-31 May 2013, pp. 6029–6033.
- [49] L. Xiao, "Dual averaging methods for regularized stochastic learning and online optimization," *J. Mach. Learn. Res.*, vol. 11, pp. 2543–2596, Oct. 2010.
- [50] J. Mairal, "Stochastic Majorization-Minimization algorithms for large-scale optimization," in *Proc. Adv. Conf. Neur. Inform. Proc. Syst.*, Lake Tahoe, Nevada, Dec. 5-8 2013, pp. x–x+8.
- [51] D. Carlson, Y.-P. Hsieh, E. Collins, L. Carin, and V. Cevher, "Stochastic spectral descent for discrete graphical models," *IEEE J. Sel. Top. Signal Process.*, vol. 10, no. 2, pp. 296–311, Mar. 2016.
- [52] J. R. Birge, X. Chen, L. Qi, and Z. Wei, "A stochastic Newton method for stochastic quadratic programs with recourse," Tech. Rep., 1995, <http://citeseerx.ist.psu.edu/viewdoc/summary?doi=10.1.1.49.4279>.
- [53] A. Bordes, L. Bottou, and P. Gallinari, "SGD-QN: Careful quasi-Newton stochastic gradient descent," *J. Mach. Learn. Res.*, vol. 10, pp. 1737–1754, Jul. 2009.
- [54] J. Yu, S. V. N. Vishwanathan, S. Günter, and N. N. Schraudolph, "A quasi-Newton approach to nonsmooth convex optimization problems in machine learning," *J. Mach. Learn. Res.*, vol. 11, pp. 1145–1200, Mar. 2010.

- [55] R. H. Byrd, S. L. Hansen, J. Nocedal, and Y. Singer, “A stochastic quasi-Newton method for large-scale optimization,” *SIAM J. Optim.*, vol. 26, no. 2, pp. 1008–1031, 2016.
- [56] R. M. Gower and P. Richtárik, “Randomized quasi-newton updates are linearly convergent matrix inversion algorithms,” Tech. Rep., 2016, <http://arxiv.org/abs/1602.01768>.
- [57] R. M. Gower and P. Richtárik, “Randomized iterative methods for linear systems,” *SIAM J. Matrix Anal. & Appl.*, vol. 36, no. 4, pp. 1660–1690, 2016.
- [58] H. Zou and T. Hastie, “Regularization and variable selection via the elastic net,” *J. R. Statist. Soc. B*, vol. 67, no. 2, pp. 301–320, 2005.
- [59] M. Nikolova and M. Ng, “Analysis of half-quadratic minimization methods for signal and image recovery,” *SIAM J. Sci. Comput.*, vol. 27, no. 3, pp. 937–966, 2005.
- [60] N. Bissantz, L. Dumbgen, A. Munk, and B. Stratmann, “Convergence analysis of generalized iteratively reweighted least squares algorithms on convex function spaces,” *SIAM J. Optim.*, vol. 19, no. 4, pp. 1828–1845, 2009.
- [61] M. Allain, J. Idier, and Y. Goussard, “On global and local convergence of half-quadratic algorithms,” *IEEE Trans. Image Process.*, vol. 15, no. 5, pp. 1130–1142, May 2006.
- [62] J. M. Bioucas-Dias and M. A. T. Figueiredo, “A new TwIST: two-step iterative shrinkage/thresholding algorithms for image restoration,” *IEEE Trans. Image Process.*, vol. 16, no. 12, pp. 2992–3004, Dec. 2007.
- [63] A. Beck and M. Teboulle, “A fast iterative shrinkage-thresholding algorithm for linear inverse problems,” *SIAM J. Imaging Sci.*, vol. 2, no. 1, pp. 183–202, 2009.
- [64] Y. Nesterov, “Gradient methods for minimizing composite objective function,” Tech. Rep., 2007, <http://hdl.handle.net/2078.1/5122>.
- [65] A. Miele and J. W. Cantrell, “Study on a memory gradient method for the minimization of functions,” *J. Optim. Theory Appl.*, vol. 3, no. 6, pp. 459–470, 1969.
- [66] E. Chouzenoux, F. Zolyniak, E. Gouillart, and H. Talbot, “A majorize-minimize memory gradient algorithm applied to X-ray tomography,” in *Proc. 20th IEEE Int. Conf. Image Process. (ICIP 2013)*, Melbourne, Australia, 15–18 Sep. 2013, pp. 1011–1015.
- [67] S. L. Gay and S. Tavathia, “The fast affine projection algorithm,” in *Proc. Int. Conf. Acoust., Speech Signal Process.*, Detroit, MI, May 9–12 1995, vol. 5, pp. 3023–3026.
- [68] D. G. Manolakis, V. K. Ingle, and S. M. Kogon, *Statistical and Adaptive Signal Processing*, Artech House, Norwood, MA, 2005.
- [69] H. J. Kushner and G. G. Yin, *Stochastic approximation and recursive algorithms and applications*, vol. 35 of *Stochastic Modelling and Applied Probability*, Springer-Verlag, New York, 2nd edition, 2003.
- [70] N. Frikha and S. Menozzi, “Concentration bounds for stochastic approximations,” *Electron. Commun. Probab.*, vol. 17, no. 47, pp. 1–15, 2012.
- [71] M. Fathi and N. Frikha, “Transport-entropy inequalities and deviation estimates for stochastic approximation schemes,” Tech. Rep., 2013, <https://arxiv.org/abs/1301.7740>.

- [72] E. Chouzenoux and J.-C. Pesquet, “Convergence rate analysis of the majorize-minimize subspace algorithm,” *IEEE Signal Process. Lett.*, vol. 23, no. 9, pp. 1284–1288, Sep. 2016.
- [73] A. Defazio, F. Bach, and S. Lacoste, “SAGA: A fast incremental gradient method with support for non-strongly convex composite objectives,” in *Proc. Ann. Conf. Neur. Inform. Proc. Syst.*, Montreal, Canada, Dec. 8-11 2014, pp. 1646–1654.
- [74] Y. Chen, Y. Gu, and A. O. Hero, “Regularized least-mean-square algorithms,” pp. x–x+7, 2010, <http://arxiv.org/pdf/1012.5066.pdf>.
- [75] S. Werner, J.A. Apolinario, and P. S. R. Diniz, “Set-membership proportionate affine projection algorithms,” *EURASIP Journal on Audio, Speech, and Music Processing*, vol. 2007, no. 1, pp. 034242, 2007.
- [76] J. Davidson, *Stochastic Limit Theory*, Oxford University Press, New York, 1994.
- [77] W. F. Stout, “The Hartman-Wintner law of the iterated logarithm for martingales,” *Ann. Math. Stat.*, vol. 41, no. 6, pp. 2158–2160, 1970.
- [78] T. L. Lai, “Martingales in sequential analysis and time series, 1945-1985,” *J. Electron. Hist. Probab. Stat.*, vol. 5, no. 1, pp. x–x+30, June 2009.
- [79] H. Robbins and D. Siegmund, “A convergence theorem for non negative almost supermartingales and some applications,” in *Optimizing Methods in Statistics*, J. S. Rustagi, Ed., pp. 233–257. Academic Press, New York, 1971.
- [80] A. M. Ostrowski, *Solution of Equations in Euclidean and Banach Spaces*, Academic Press, London, 1973.